

# PRISM: Private Verifiable Set Computation over Multi-Owner Outsourced Databases

## ABSTRACT

This paper proposes PRISM, a secret sharing based approach to compute private set operations (*i.e.*, intersection and union, as well as aggregates) over outsourced databases belonging to multiple owners. PRISM enables data owners to pre-load the data onto servers, exploits the additive and multiplicative properties of secret-shares to compute the above-listed operations in (at most) two rounds of communication between non-colluding servers (storing the secret-shares) and the querier, resulting in a very efficient implementation. PRISM also supports result verification techniques for each operation to detect malicious adversaries. Experimental results show that PRISM scales both in terms of the number of data owners and database sizes, to which prior approaches do not scale.

## 1 INTRODUCTION

With the advent of cloud computing, database-as-a-service (DaS) [31] has gained significant attention. Traditionally, the DaS problem focused on a single database (DB) owner, submitting suitably encrypted data to the cloud over which DB owner (or one of its clients) can execute queries. A more general use-case is one in which there are multiple datasets, each owned by a different owner. Data owners do not trust each other, but wish to execute queries over common attributes of the dataset. The query execution must not reveal the content of the database belonging to one DB owner to others, except for the leakage that may occur from the answer to the query. The most common form of such queries is the *private set intersection* (PSI) [24, 34, 39, 43, 56, 57, 59]. An example use-case of PSI include *syndromic surveillance*, wherein organizations, such as pharmacies and/or hospitals, share information (*e.g.*, a sudden increase in sales of specific drugs such as analgesics or anti-allergy medicine, tele-health calls, and school absenteeism requests) to enable early detection of community-wide outbreaks of diseases. PSI is also a building block for performing joins across private databases — it essentially corresponds to performing a semi-join operation on the join attribute [40].

Private set computations over datasets owned by different DB owners/organizations can, in general, be implemented using secure multiparty computation (SMC) [28, 49, 68], a well-known cryptographic technique that has been prevalent since more than three decades. SMC allows DB owners to securely execute any function over their datasets without revealing their data to other DB owners. However, SMC can be very slow, often by order of magnitude [45]. As a result, techniques that can be used to more efficiently compute private set operations have been developed; particularly, in the context of PSI [24, 34, 39, 56, 59] and *private set union* (PSU) [19, 42]. PSU refers to privately computing the union of all databases. Several approaches using homomorphic encryption [14], polynomial evaluation [24], garbled-circuit techniques [34], hashing [23, 56, 64], hashing and oblivious pseudo-random functions (OPRF) [44], Bloom-filter [52], and oblivious transfer [55, 58] have been proposed to implement private set operations.

Recent work on private set operations has also explored performing aggregation on the result of PSI operations. For instance, [36] studied the problem of private set intersection sum (PSI Sum), motivated by

the internet advertising use-case, where a party maintains information about which customer clicked on specific advertisements during their web session, while another has a list of transactions about items listed in the advertisements that resulted in a purchase by the customers. Both parties might wish to securely learn the total sales that can be attributed to customers clicking on advertisements, while neither would like their databases to be revealed to the other for reasons including fair/competitive business strategies.

Existing approaches on private set computation (including recent work on aggregation) are limited in several ways:

- Work on PSI or PSU has largely focused on the case of two DB owners, with some exceptions that address more than two DB owners scenarios, *e.g.*, [15, 24, 33, 35, 42, 45, 69]. There are several interesting use-cases, where one may wish to compute PSI over multiple datasets. For instance, in the syndromic surveillance example listed above, one may wish to compute intersection amongst several independently owned databases. Generalizing existing two-party PSI or PSU approaches to the case of multiple DB owners results in significant overhead [45]. For instance, [2], which is designed for two DB owners, incurs  $(nm)^2$  communication cost, when extended to  $m > 2$  DB owners, where  $n$  is the dataset size. Even recent work supporting multiple DB owners incurred significant computational overhead; *e.g.*, [35] took  $\approx 12$  seconds for PSI over 24 DB owners having 1024 values.
- Techniques to privately compute aggregation over set operations have not been studied systematically. In the database literature, aggregation functions [51] are typically classified as: *summary* aggregations (such as count, sum, and average) or *exemplary* aggregations (such as minimum, maximum, and median). Existing literature has only considered the problem of PSI Sum [36] and cardinality determination, *i.e.*, the size of the intersection or union [19, 22]. Techniques for exemplary aggregations (and even for summary aggregations) that may compute over multiple attributes have not been explored.
- Many of the existing solutions do not deal with a large amount of data, due to either inefficient cryptographic techniques or multiple communication rounds amongst DB owners. For instance, recent work [45, 46, 69] dealt with data that is limited to sets of size less than or equal to  $\approx 1M$  in size.

This paper introduces PRISM — a novel approach for computing collaboratively over multiple databases. PRISM is designed for both PSI and PSU, and it supports both summary, as well as, exemplar aggregations. Unlike existing SMC techniques (wherein DB owners compute operations privately through a sequence of communication rounds), in PRISM, DB owners outsource their data in secret-shared form to multiple *non-communicating public servers*. As will become clear, PRISM exploits the homomorphic nature of secret-shares (both additive and multiplicative) to enable servers to compute private set operations independently (to a large degree). These results are then returned to DB owners to compute the final results. In PRISM, any operator requires at most two communication rounds between DB owners and servers, where the first round finds tuples that are in the intersection or union of the set, and the second round computes the aggregation function over the objects in the intersection/union.

	Name	Age	Disease	Cost
$\tau_1$	John	4	Cancer	100
$\tau_2$	Adam	6	Cancer	200
$\tau_3$	Mike	2	Heart	300

**Table 1: Hospital 1.**

	Name	Age	Disease	Cost
$v_1$	John	8	Cancer	100
$v_2$	Adam	5	Fever	70
$v_3$	Bob	4	Fever	50

**Table 2: Hospital 2.**

Note:  $\tau_i$ ,  $v_i$ , and  $\rho_i$  denote the  $i^{th}$  tuples of tables.

	Name	Age	Disease	Cost
$\rho_1$	Carl	8	Cancer	300
$\rho_2$	John	4	Cancer	700
$\rho_3$	Lisa	5	Heart	500

**Table 3: Hospital 3.**

By using public servers for computation over secret-shared data, PRISM achieves the identical security guarantees as existing SMC systems (e.g., Sharemind [7], Jana [4], and Conclave [65]). The key advantage of PRISM is that by outsourcing data in secret shared form and exploiting homomorphic properties, PRISM does not require communication among server before/during/after the computation, which allows PRISM to perform efficiently even for large data sizes and for a large number of DB owners (as we will show in experiment section). Since PRISM uses the public servers, which may act maliciously, PRISM supports oblivious result verification methods.

In summary, PRISM offers the following benefits: (i) *Information-theoretical security*: It achieves information-theoretical security at the servers and prevents them to learn anything from input/output/access-patterns/output-size. (ii) *No communication among servers*: It does not require any communication among servers, unlike SMC based solutions. (iii) *No trusted entity*: It does not require any trusted entity that performs the computation on the cleartext data, unlike the recent SMC system Conclave [65]. (iv) *Several DB owners and large-sized dataset*: It deals with several DB owners having a large-size dataset. **Full version.** Due to space restriction, we omit the following from this version and provide in the full version: PSU result verification, PSU-count, PSU-sum, PSU-average, an example for PSU, maximum verification, top-k, and median verification.

## 2 PRIVATE SET OPERATIONS

We, first, define the set of operations supported by PRISM. Let  $DB_1, \dots, DB_m$  be  $m > 2$  independent DBs owned by  $m$  DB owners  $\mathcal{DB}_1, \dots, \mathcal{DB}_m$ . We assume that each DB owner is (partially) aware of the schema of data stored at other DB owners. Particularly, DB owners have knowledge of the attribute(s) of the data stored at other DB owners on which the set-based operations (i.e., intersection or union) can be performed. Also, DB owners know about the attributes on which aggregation functions (e.g., sum, min, max) be supported. Such an assumption is needed to ensure that PSI/PSU and aggregation queries are well defined. Other than the above requirement, the schema of data at different databases may be different.

Now, we define the private set operations supported by PRISM formally and their corresponding privacy requirements (corresponding SQL statements are shown in Table 4). In defining the operators (and in the rest of the paper), we will use the example tables shown in Tables 1, 2, and 3 that are owned by three different DB owners (in our case, hospitals).

**Private set intersection (PSI) (§5).** PSI finds the common values among  $m$  DB owners for a specific attribute  $A_c$ , i.e.,  $DB_1.A_c \cap \dots \cap DB_m.A_c$ . For example, PSI over disease column of Tables 1, 2, and 3 will return {Cancer} as a common disease treated by all hospitals. Note that a hospital computing PSI on disease should not gain any information about other possible disease values (except for the result of the PSI) associated with other hospitals.

**Private set union (PSU) (§7).** PSU finds the union of values among  $m$  DB owners for a specific attribute  $A_c$ , i.e.,  $DB_1.A_c \cup \dots \cup DB_m.A_c$ . For example, PSU over disease column returns {Cancer, Fever, Heart} as diseases treated by all hospitals. Again, a hospital computing PSU over other hospitals must not gain information about the specific diseases treated by other hospitals, or how many hospitals treat which disease.

**Aggregation over private set operators (§6.)** Aggregation  $A_c \mathcal{G}_\theta(A_x)$  computes the aggregation function  $\theta$  on the attribute  $A_x$  ( $A_c \neq A_x$ ) for the groups corresponding to the output of set-based operations (PSI or PSU) on attribute  $A_c$ . For example, the aggregation function *sum* on cost attribute corresponding to PSI over disease attribute (i.e.,  $\text{disease} \mathcal{G}_{\text{sum}}(\text{cost})$ ) returns a tuple {Cancer, 1400}. The same aggregation function over PSU will return {(Cancer, 1400), (Fever, 120), (Heart, 800)}. Likewise, the output of aggregation  $\text{disease} \mathcal{G}_{\text{max}}(\text{age})$  over PSI would return {Cancer, 8}, while the same over PSU would return {(Cancer, 8), (Fever, 5), (Heart, 5)}. Note that the count operation does not require specifying an aggregation attribute  $A_x$  and can be computed over the attribute(s) associated with PSI or PSU. For example, count over PSI (PSU) on disease column will return 1 (3) respectively. From the perspective of privacy requirement, in case of PSI on disease column, a hospital executing an aggregation query (maximum of age or sum of cost) should only gain information about the answer, i.e., *elements in the PSI and the corresponding aggregate value*. It should not gain information about other diseases that are not in the intersection. Likewise, for PSU, the hospital will gain information about *all elements in the union and their corresponding aggregate values, but will not gain any specific information about which database contains which disease values, or the number of databases with a specific disease*.

## 3 PRELIMINARY

This section describes the cryptographic concepts that serve as building blocks for PRISM, provides an overview of PRISM approach, and discusses its security properties.

### 3.1 Building Blocks

PRISM is based on additive secret-sharing (SS), Shamir's secret-sharing (SSS), cyclic group, and pseudorandom number generator. We provide an overview of these techniques, below.

**Additive Secret-Sharing (SS).** Additive SS is the simplest type of the SS. Let  $\delta$  be a prime number. Let  $\mathbb{G}_\delta$  be an Abelian group under modulo addition  $\delta$  operation. All additive shares are defined over  $\mathbb{G}_\delta$ . In particular, the DB owner creates  $c$  shares  $A(s)^1, A(s)^2, \dots, A(s)^c$  over  $\mathbb{G}_\delta$  of a secret, say  $s$ , such that  $s = A(s)^1 + A(s)^2 + \dots + A(s)^c$ . The DB owner sends share  $A(s)^i$  to the  $i^{th}$  server (belonging to a set of  $c$  non-communicating servers). These servers cannot know the secret  $s$  until they collect all  $c$  shares. To reconstruct  $s$ , the DB owner collects all the shares and adds them. Additive SS allows *additive homomorphism*. Thus, servers holding shares of different secrets can locally compute the sum of those shares. Let  $A(x)^i$  and  $A(y)^i$  be additive shares of two secrets  $x$  and  $y$ , respectively, at a server  $i$ , then the server  $i$  can compute  $A(x)^i + A(y)^i$  that enable DB owner to know the result of  $x+y$ . The precondition of *additive homomorphism* is that the sum of shares should be less than  $\delta$ .

**Example.** Let  $\mathbb{G}_5 = \{0, 1, 2, 3, 4\}$  be an Abelian group under the addition modulo 5. Let 4 be a secret. The DB owner may create two shares, such as 3 and 1 (since  $4 = (3+1) \bmod 5$ ), and sends them to two servers.

**Shamir's Secret-Sharing (SSS) [62].** In SSS [62], the DB owner randomly selects a polynomial of degree  $c'$  with  $c'$  random coefficients, i.e.,  $f(x) = a_0 + a_1x + a_2x^2 + \dots + a_{c'}x^{c'}$ , where  $f(x) \in \mathbb{F}_p[x]$ ,  $p$  is a prime number,  $\mathbb{F}_p$  is a finite field of order  $p$ ,  $a_0 = s$ , and  $a_i \in \mathbb{N}$ .

PSI	SELECT $A_c$ FROM $db_1$ INTERSECT ... INTERSECT SELECT $A_c$ FROM $db_m$
PSU	SELECT $A_c$ FROM $db_1$ UNION ... UNION SELECT $A_c$ FROM $db_m$
PSI count	SELECT COUNT( $A_c$ ) FROM $db_1$ INTERSECT ... INTERSECT SELECT $A_c$ FROM $db_m$
PSI $\theta$ $\theta \in \{AVG, SUM, MAX, MIN, Median\}$	CREATE VIEW $CommonA_c$ as SELECT $A_c$ FROM $db_1$ INTERSECT ... INTERSECT SELECT $A_c$ FROM $db_m$ SELECT $A_c, \theta(A_x)$ FROM (SELECT $A_x, A_c$ FROM $db_1, CommonA_c$ WHERE $db_1.A_c = CommonA_c.A_c$ UNION ALL ... UNION ALL SELECT $A_x, A_c$ FROM $db_m, CommonA_c$ WHERE $db_m.A_c = CommonA_c.A_c$ ) as inner_relation Group By $A_c$

**Table 4: SQL syntax of operations supported by PRISM.**

( $1 \leq i \leq c'$ ). The DB owner distributes the secret  $s$  into  $c$  shares by computing  $f(x)$  ( $x = 1, 2, \dots, c$ ) and sends an  $i^{th}$  share to the  $i^{th}$  server (belonging to a set of  $c$  non-communicating servers). The secret can be reconstructed using any  $c' + 1$  shares using Lagrange interpolation [17].

SSS, also, allows *additive homomorphism*, i.e., if  $S(x)^i$  and  $S(y)^i$  are SSS of two secrets  $x$  and  $y$ , respectively, at a server  $i$ , then the server  $i$  can compute  $S(x)^i + S(y)^i$ , which will result in  $x + y$  at DB owner.

**Cyclic group under modulo multiplication.** Let  $\eta$  be a prime number. A group  $\mathbb{G}$  is called a cyclic group, if there exists an element  $g \in \mathbb{G}$ , such that all  $x \in \mathbb{G}$  can be derived as  $x = (g^i)$  (where  $i$  is an integer number  $\mathbb{Z}$ ) under modulo multiplicative  $\eta$  operation. The element  $g$  is called a generator of the cyclic group, and the number of elements in  $\mathbb{G}$  is called the *order* of  $\mathbb{G}$ . Based on each element  $x$  of a cyclic group, we can form a cyclic subgroup by executing  $x^i \bmod \eta$ . *Example.*  $g = 2$  is a generator of a cyclic group under multiplication modulo  $\eta = 11$  for the following group:  $\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$ . Note that the group elements are derived by  $2^i \bmod 11$ . By taking 5 of this cyclic group, we form the following cyclic subgroup  $\{1, 3, 4, 5, 9\}$ , under multiplication modulo  $\eta = 11$ , by  $5^i \bmod 11$ .

**Permutation function  $\mathcal{PF}$ .** Let  $A$  be a set. A permutation function  $\mathcal{PF}$  is a bijective function that maps a permutation of  $A$  to another permutation of  $A$ , i.e.,  $\mathcal{PF}: A \rightarrow A$ .

**Pseudorandom number generator  $\mathcal{PRG}$ .** A pseudorandom number generator is a deterministic and efficient algorithm that generates a pseudorandom number sequence based on an input seed [6, 27].

### 3.2 Entities and Trust Assumption

PRISM assumes the following four entities:

- (1) The  $m$  **database (DB) owners** (or users), who wish to execute computation on their joint datasets. We assume that each DB owner is trusted and does not act maliciously.
- (2) A set of  $c \geq 2$  **servers** that store the secret-shared data outsourced by DB owners and execute the requested computation from authenticated DB owners. The data transmission between a DB owner and a server takes place in encrypted form or by using anonymous routing [29] to prevent the locations of servers and the shares from an adversary, eavesdropping on the communication channel between DB owners and servers.
- We assume that servers do not maliciously communicate (i.e., non-communicating servers) with each other in violation of PRISM protocols. Unlike other MPC mechanisms [7], (as will be clear soon), PRISM protocols do not require the servers to communicate before/during/after the execution of the query. The security of secret-sharing techniques requires that out of the  $c$  servers, no more than  $c' < c$  communicate maliciously or collude with each other, where  $c'$  is a minority of servers (i.e., less than half of  $c$ ). Thus, we assume that a majority of servers do not collude and communicate with each other, and hence, a legal secret value cannot be generated/inserted/updated/deleted at the majority of the servers. Also, note that the collusion of servers in violation of the protocol is a general requirement for secret-sharing based protocols, and a similar assumption is made by many prior work [7, 16, 62, 67]. This assumption is based on factors such as economic incentivization (violation is against their economic interest), law (illegal to collude),

and jurisdictional boundaries. Such servers can be selected on different clouds, which make the assumption more realistic.

For the purpose of simplicity, we assume that none of the servers collude with each other – that is they not communicate directly. Thus, to reconstruct the original secret value from the shares, *two additive shares* suffice. In the case of PSI sum (as will be clear in §6.1), we need to multiply two shares (each of degree one) and that increases the degree of the polynomial to two. To reconstruct the secret value of degree two, we need at least three multiplicative (Shamir’s secret) shares.

While we assume that servers do not collude, we will consider two types of adversarial models for the servers in the context of the computation that they perform: (i) **Honest-but-curious (HBC)** servers that correctly compute the assigned task without tampering with data or hiding answers. It may, however, exploit side information (e.g., the internal state of the server, query execution, background knowledge, and output size) to gain as much information as possible about the stored data/computation/results. The HBC adversarial model is considered widely in many cryptographic algorithms and in DaS model [12, 31, 66]. (ii) **Malicious** adversarial servers that can delete/insert tuples from the relation, and hence, is a stronger adversarial model than HBC.

(3) An **initiator or oracle**, who knows  $m$  DB owners and servers. Before data outsourcing by DB owners, the initiator informs the identity of servers to DB owners and vice versa. Also, the initiator informs the desired parameters (e.g., a hash function, parameters related to Abelian and cyclic groups,  $\mathcal{PF}$ , and  $\mathcal{PRG}$ ) to servers and DB owners. The initiator is an entity trusted by all other entities and plays a role similar to the trusted certificate authority in the public-key infrastructure. The initiator never knows the data/results, since it does not store any data, or data/results are not provided to servers via the initiator. The role of the initiator has also been considered in existing PSI work [60, 70].

(4) An **announcer  $S_a$**  that participates only in maximum, minimum, and median queries to announce the results.  $S_a$  communicates (not maliciously) with servers and the initiator (and not with DB owners).

### 3.3 PRISM Overview

Let us first understand the working of PRISM at the high-level. PRISM contains four phases (see Figure 1), as follows:

**PHASE 0: Initialization.** The initiator sends desired parameters (see details in §4) related to additive SS, SSS, cyclic group,  $\mathcal{PF}$ , and  $\mathcal{PRG}$  to all entities and informs them about the identity of others from/to whom they will receive/send the data.

**PHASE 1: Data Outsourcing by DB owners.** DB owners create additive SS or SSS of their data, by following the methods given in §5 for PSI and PSU, §6.1 for PSI/PSU-sum, and §6.3 for PSI/PSU-maximum/minimum. Then, they outsource their secret-shared data to non-communicating servers. Note that for the purpose of explanations, we will write the data outsourcing method with query execution.

**PHASE 2: Query Generation by the DB owner.** A DB owner who wishes to execute SMC over datasets of different DB owners, sends the query to the servers. For generating secret-shared queries for PSI, PSU, count, sum, maximum, and for their verification, the DB owner follows the method given in §5, 6.



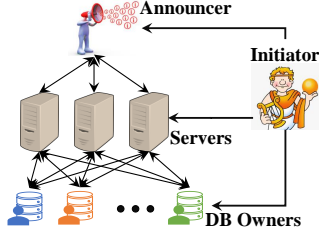


Figure 1: PRISM model.

**PHASE 3: Query Processing.** The servers process an input query and respective verification method in an oblivious manner, such that neither the query nor the results satisfying the query/verification are revealed to the adversary. Finally, servers transfer their outputs to DB owners.

**PHASE 4: Final processing at the DB owners.** The DB owner either adds the additive shares or performs Lagrange interpolation on SSS to obtain the answer to the query.

### 3.4 Security Property

As mentioned in the adversarial setting in §3.2, an adversarial server wishes to learn the (entire/partial) input and output data, while a DB owner may wish to know the data of other DB owners. Hence, a secure algorithm must prevent an adversary to learn the data (i) from the ciphertext representation of the data, (ii) from query execution due to access-patterns (i.e., the adversary can learn the physical locations of tuples that are accessed to answer the query), and (iii) from the size of the output (i.e., the adversary can learn the number of tuples satisfy the query). The attacks on a dataset based on access-patterns and output-size are discussed in [13, 37]. In order to prevent these attacks, our security properties are identical to the standard security definition as in [11, 12, 25]. An algorithm is *privacy-preserving* if it maintains DB owners' privacy, data/computation privacy from the servers, and performs identical operations regardless of the inputs.

**Privacy from servers** requires that datasets of DB owners must be hidden from the server, before/during/after any computation. In PSI/PSU, servers must not know whether a value is common or not, the number of DB owners having a particular value in the result set. In the case of aggregation operations, the output of aggregation over an attribute  $A_x$  corresponding to the attributes  $A_c$  involved in PSI or PSU should not be revealed to servers. Additionally, in the case of maximum/median/minimum query, servers must not know the maximum/minimum value and the identity of the DB owner who possesses such values. Further, the protocol must ensure that the server's behavior in reading/sending the data must be identical for a particular type of query (e.g., PSI or PSU), thereby the server should not learn anything from query execution (i.e., hiding access-patterns and output-sizes).

**DB owner privacy** requires that the DB owners must not learn anything other than their datasets and the final output of the computation. For example, in PSI/PSU queries, DB owners must only learn the intersection/union set, and they must not learn the number of DB owners that does not contain a particular value in their datasets. Similarly, in the case of aggregation operations, DB owners must only learn the output of aggregation operation, not the individual values on which aggregation was performed.

**Properties of verification.** A verification method must be oblivious and find any misbehavior of servers when computing a query. We follow the verification properties from [38] that the verification

Notation	Meaning
$A_c$	An attribute on which PSI/PSU executes
$\text{Dom}(A_c)$	domain of attribute $A_c$
$A_x$	An attribute used in aggregation
$m$	Number of DB owners
$\mathcal{DB}_i$	The $i^{\text{th}}$ DB owner
$A(m)^\phi$	$\phi^{\text{th}}$ additive share of the number $m$
$\delta$	A prime number defining modulo addition $\delta$ operation
$\mathbb{G}_\delta$	Abelian group under modulo addition $\delta$ operation
$g$	Generator of a cyclic group
$\eta$	A prime number defining modulo multiplicative $\eta$ operation
$\eta'$	$\eta \times \alpha$ , where $\alpha > 1$
$\mathcal{PRG}$	Pseudo-random number generator
$S_a$	The announcer
$S_\phi$	A server $\phi$ , where $\phi \in \{1, 2, 3\}$
$\mathcal{PF}_{s1}, \mathcal{PF}_{s2}$	The permutation functions known to servers
$\mathcal{PF}_{db1}, \mathcal{PF}_{db2}$	The permutation functions known to all DB owners.
$\mathcal{PF}_i$	A permutation function known to the initiator
$\mathcal{PF}$	A permutation function known to both servers and DB owners
$\chi = \{x_1, x_2, \dots, x_b\}$	A hash table at the $j^{\text{th}}$ DB owner of length $b =  \text{Dom}(A_c) $
$A(x_j)^\phi$	The $\phi^{\text{th}}$ additive share of an $i^{\text{th}}$ element of $\chi_j$ of $\mathcal{DB}_j$
$\bar{x}_j$	The complement value of $x_j$
$\bar{\chi}_j$	A table having the complement values of $\chi_j$
$A(\chi_j)^\phi$	The $\phi^{\text{th}}$ additive share of the table $\chi_j$
$A(\bar{\chi}_j)^\phi$	The $\phi^{\text{th}}$ additive share of the table $\bar{\chi}_j$
$\text{output}_i^{S_\phi}$	The output computed for the $i$ -th value at server $S_\phi$
$\ominus$	The modular subtraction operation
$\oplus$	The modular addition operation

Table 5: Frequently used notations in the paper.

method cannot be refuted by the majority of the malicious servers and should not leak any additional information.

## 4 ASSUMPTIONS & PARAMETERS

Different entities in PRISM protocols are aware of the following parameters to execute the desired task:

**Parameters known to the initiator.** The initiator knows all parameters used in PRISM and distributes them to different entities (only once) as they join in PRISM protocols. Note that the initiator can select these parameters (such as  $\eta$ ,  $\delta$ ) to be large to support increasing DB owners over time without updating parameters. Thus, when new DB owners join, the initiator simply needs to inform DB owners/servers about the increase in the number of DB owners in the system, but does not need to change *all* parameters.

Additionally, the initiator does the following: (i) Selects a polynomial ( $\mathcal{F}(x) = a_{m+1}x^{m+1} + a_mx^m + \dots + a_1x + a_0$ , where  $a_i > 0$ ) of degree more than  $m$ , where  $m$  is the number of DB owners, and sends the polynomial to all DB owners. This polynomial will be used during the maximum computation. Importantly, this polynomial  $\mathcal{F}(x)$  generates values at different DB owners in an order-preserving manner, as will be clear in §6.3, and the degree of the polynomial must be more than  $m$  to prevent an entity, who has  $m$  different values generated using this polynomial, to reconstruct the secret value (a condition similar to SSS); and beyond  $m + 1$ , the degree of the polynomial does not impact the security, in this case. (ii) Generates a permutation function  $\mathcal{PF}_i$ , and produces four different permutation functions that satisfy Equation 1:

$$\mathcal{PF}_{s1} \odot \mathcal{PF}_{db1} = \mathcal{PF}_{s2} \odot \mathcal{PF}_{db2} = \mathcal{PF}_i \quad (1)$$

Here, the symbol  $\odot$  represents composition of the permutations, and these functions can be selected over a permutation group. The initiator provides  $\mathcal{PF}_{s1}$  and  $\mathcal{PF}_{s2}$  to all servers and  $\mathcal{PF}_{db1}$  and  $\mathcal{PF}_{db2}$  to all DB owners.

**Parameters known to the announcer.** The announcer  $S_a$  knows  $\delta$ , a prime number used to define modulo addition for an Abelian group (§3.1). The announcer helps in maximum and median algorithms.

**Parameters known to DB owners.** All DB owners know the following parameters: (i)  $m$ , i.e., the number of DB owners. (ii)  $\delta > m$ , (iii)  $\eta$ , where  $\eta$  is a prime number used to define modular multiplication for a cyclic group (§3.1). Note that DB owners do not know the generator  $g$  of the cyclic group. (iv) A common hash function. (v) The domain of the attribute  $A_c$  on which they want to execute PSI/PSU. Note that knowing the domain of the attribute  $A_c$  does not reveal that which of the DB owner has a value of the domain or not. (Such an assumption is also considered in prior work [34].) (vi) Two permutation functions  $\mathcal{PF}_{db1}$  and  $\mathcal{PF}_{db2}$ . (vii) The polynomial  $\mathcal{F}(x)$  given by the initiator. (viii) A permutation function  $\mathcal{PF}$ , and the same permutation function will also known to servers.

PSI, PSU, sum, average, count algorithms are based on the assumptions 1-5. PSI verification, sum verification, count, and count verification algorithms are based on the assumptions 1-6. Maximum, maximum verification, and median algorithms are based on the assumptions 1-8.

Further, we assume that any DB owner or the initiator provides additive shares of  $m$  to servers for executing PSI, and the DB owners have only positive integers to compute the maximum. Since the current PSI maximum method uses modular operations (as will be clear in §6.3), we cannot handle floating point values directly. Nevertheless, we can find the maximum for a large class of practical situations, where the precision of decimal is limited, say  $k > 0$  digits by simply multiplying each number by  $10^k$  and using the current PSI maximum algorithm. For example, we can find the maximum over  $\{0.5, 8.2, 8.02\}$  by computing the maximum over  $\{50, 820, 802\}$ . Designing a more general solution that does not require limited precision is non-trivial.

**Parameters known in servers.** Servers know the following parameters: (i)  $m$ ,  $\delta > m$ , the generator  $g$  of the cyclic (sub)group of order  $\delta$  and  $\eta' = \alpha \times \eta$  and  $\alpha > 1$ . Also, based on the group theory,  $\eta - 1$  should be divisible by  $\delta$ . Note that servers do not know  $\eta$ . (ii) A permutation function  $\mathcal{PF}$ , and recall that the same permutation function is also known to DB owners. (iii) Two permutation functions  $\mathcal{PF}_{s1}$  and  $\mathcal{PF}_{s2}$ . (iv) A common pseudo-random number generator  $\mathcal{PRG}$  that generates random numbers between 1 and  $\delta - 1$ . Note that  $\mathcal{PRG}$  is unknown to DB owners. PSI, sum, and average algorithms are based on the assumptions 1. Maximum, maximum verification, and median algorithms are based on the assumptions 1, 2. Count and its verification are based on the assumptions 1, 3. PSU and its verification are based on the assumptions 1, 4.

## 5 PRIVATE SET INTERSECTION (PSI) QUERY

This section, first, develops a method for finding PSI among  $m > 2$  different DB owners on an attribute  $A_c$  (which is assumed to exist at all DB owners, §5.1) and presents a result verification method (§5.2). Later in §6.6, we present a method to execute PSI over multiple attributes and a method to reduce the communication cost of PSI.

### 5.1 PSI Query Execution

**High-level idea.** Each of  $m > 2$  DB owners uses a publicly known hash function to map distinct values of  $A_c$  attribute in a table of cells at most  $|\text{Dom}(A_c)|$ , where  $|\text{Dom}(A_c)|$  refers to the size of the domain of  $A_c$ . Thus, if a value  $a_j \in A_c$  exists at any DB owner, all DB owners must map  $a_j$  to an identical cell of the table. Then, all values of the table are outsourced in the form of additive shares to two non-communicating servers  $\mathcal{S}_\phi$ ,  $\phi \in \{1, 2\}$ , that obviously find the common items/intersection and return shared output vector (of

the same length as the length of the received shares from DB owners). Finally, each DB owner adds the results to know the final answer.

**Construction.** We create the following construction over the elements of a group under addition and the elements of a cyclic group under multiplication. Note that we can select any cyclic group such that  $\eta > m$ .

$$(x+y) \bmod \delta = 0, (g^x \times g^y) \bmod \eta = 1 \quad (2)$$

Based on this construction, below, we explain PSI finding algorithm:

**STEP 1: DB owners.** Each DB owner finds distinct values in an attribute ( $A_c$ , which exists at all DB owners, as per our assumption given in §4) and executes the hash function on each value  $a_i$  to create a table  $\chi = \{x_1, x_2, \dots, x_b\}$  of length  $b = |\text{Dom}(A_c)|$ . The hash function maps the value  $a_i \in A_c$  to one of the cells of  $\chi$ , such that the cell of  $\chi$  corresponding to the value  $a_i$  holds 1; otherwise 0. It is important that each cell must contain only a single one corresponding to the unique value of the attribute  $A_c$ , and note that if a value  $a_i \in A_c$  exists at any DB owner, then one corresponding to  $a_i$  is placed at an identical cell of  $\chi$  at the DB owner. The table at  $\mathcal{DB}_j$  is denoted by  $\chi_j$ . Finally,  $\mathcal{DB}_j$  creates additive secret-shares of each value of  $\chi_j$  (i.e., additive secret-shares of either one or zero) and outsources the  $\phi^{th}$ ,  $\phi \in \{1, 2\}$ , share to the server  $\mathcal{S}_\phi$ . We use the notation  $A(x_i)_j^\phi$  to refer to  $\phi^{th}$  additive share of an  $i^{th}$  element of  $\chi_j$  of  $\mathcal{DB}_j$ . Recall that before the computation starts, the initiator informs the locations of servers to DB owners and vice versa (§3.2).

**STEP 2: Servers.** Each server  $\mathcal{S}_\phi$  ( $\phi \in \{1, 2\}$ ) holds the  $\phi^{th}$  additive share of the table  $\chi$  (denoted by  $A(\chi)_j^\phi$ ) of  $j^{th}$  ( $1 \leq j \leq m$ ) DB owners and executes Equation 3:

$$output_i^{\mathcal{S}_\phi} \leftarrow g^{((\oplus_{j=1}^m A(x_i)_j^\phi) \ominus A(m)^\phi) \bmod \eta', (1 \leq i \leq b)} \quad (3)$$

where  $\oplus$  and  $\ominus$  show the modular addition and modular subtraction operations, respectively. We used the symbols  $\oplus$  and  $\ominus$  to distinguish them from the normal addition and subtraction. Particularly, each server  $\mathcal{S}_\phi$  performs the following operations: (i) modular addition (under  $\delta$ ) of the  $i^{th}$  additive secret-shares from all  $m$  DB owners, (ii) modular subtraction (under  $\delta$ ) of the result of the previous step from the additive share of  $m$  (i.e.,  $A(m)^\phi$ ), (iii) exponentiation by  $g$  to the power the result of the previous step and modulo by  $\eta'$ , and (iv) sends all the computed  $b$  results to the  $m$  DB owners.

**STEP 3: DB owners.** From two servers, DB owners receive two vectors, each of length  $b$ , and perform modular multiplication (under  $\eta$ ) of outputs  $output_i^{\mathcal{S}_1}$  and  $output_i^{\mathcal{S}_2}$ , where  $1 \leq i \leq b$ , i.e.,

$$fop_i \leftarrow (output_i^{\mathcal{S}_1} \times output_i^{\mathcal{S}_2}) \bmod \eta \quad (4)$$

This step results in an output array of  $b$  elements, which may contain any value. However, if an  $i^{th}$  item of  $\chi_j$  exists at all DB owners, then  $fop_i$  must be one, since  $\mathcal{S}_\phi$  have added additive shares of  $m$  ones at the  $i^{th}$  element and subtracted from additive share of  $m$  that results in  $(g^0 \bmod \eta') \bmod \eta = 1$  at DB owner. Please see the correctness argument below after the example.

**Example 5.1.** Assume three DB owners:  $\mathcal{DB}_1$ ,  $\mathcal{DB}_2$ , and  $\mathcal{DB}_3$ ; see Tables 1, 2, and 3. For answering a query to find the common disease that is treated by each hospital, DB owners create their tables  $\chi$  as shown in the first column of Tables 6, 7, and 8. For example, in Table 7,  $\langle 1, 1, 0 \rangle$  corresponds to cancer, fever, and heart diseases, where 1 means that the disease is treated by the hospital. We select  $\delta = 5$ ,  $\eta = 11$ , and  $\eta' = 143$ . Hence, the Abelian group under modulo

Value	Share 1	Share 2	Value	Share 1	Share 2	Value	Share 1	Share 2
1	4	-3	1	3	-2	1	2	-1
0	2	-2	1	4	-3	0	3	-3
1	3	-2	0	3	-3	1	4	-3

Table 6:  $\mathcal{DB}_1$ .Table 7:  $\mathcal{DB}_2$ .Table 8:  $\mathcal{DB}_3$ .

addition contains  $\{0,1,2,3,4\}$ , and the cyclic (sub)group (with  $g=3$ ) under modulo multiplication contains  $\{1,3,4,5,9\}$ . Assume additive shares of  $m=3=(1+2) \bmod 5$ .

**Step 1: DB Owners.** DB owners generate additive shares as shown in the second and third columns of Tables 6, 7, and 8, and outsource all values of the second and third columns to  $\mathcal{S}_1$  and  $\mathcal{S}_2$ , respectively.

**Step 2: Servers.** The server  $\mathcal{S}_1$  will return the three values 27, 27, 81, by executing the following computation, to all three DB owners:

$$\begin{aligned} &3(((4+3+2) \bmod 5)-1) \bmod 5 \bmod 143=27 \\ &3(((2+4+3) \bmod 5)-1) \bmod 5 \bmod 143=27 \\ &3(((3+3+4) \bmod 5)-1) \bmod 5 \bmod 143=81 \end{aligned}$$

The server  $\mathcal{S}_2$  will return values 9, 1, and 1 to all three DB owners:

$$\begin{aligned} &3((( -3-2-1) \bmod 5)-2) \bmod 5 \bmod 143=9 \\ &3((( -2-3-3) \bmod 5)-2) \bmod 5 \bmod 143=1 \\ &3((( -2-3-3) \bmod 5)-2) \bmod 5 \bmod 143=1 \end{aligned}$$

**Step 3: DB owners.** The DB owner obtains a vector  $\langle 1, 5, 4 \rangle$ , by executing the following computation (see below). From the vector  $\langle 1, 5, 4 \rangle$ , DB owners learn that cancer is a common disease treated by all three hospitals. However, the DB owner does not learn anything more than this; note that in the output vector, the values 5 and 4 correspond to zero. For instance,  $\mathcal{DB}_1$ , i.e., hospital 1, cannot learn whether fever and heart diseases are treated by hospital 2, 3, or not.

$$\begin{aligned} &(27 \times 9) \bmod 11 = 1 \\ &(27 \times 1) \bmod 11 = 5 \\ &(81 \times 1) \bmod 11 = 4 \blacksquare \end{aligned}$$

**Correctness.** When we plug Equation 3 into Equation 4, we obtain:

$$\begin{aligned} fop_i &= (g^{(\oplus_{j=1}^{j=m} A(x_i)_j^1) \ominus A(m)^1} \times g^{(\oplus_{j=1}^{j=m} A(x_i)_j^2) \ominus A(m)^2} \bmod \eta') \bmod \eta \\ &= (g^{(\oplus_{j=1}^{j=m} (x_i)_j - m)} \bmod \eta') \bmod \eta \end{aligned}$$

We utilize the modular identity, i.e.,  $(x \bmod a\eta) \bmod \eta = x \bmod \eta$ ;

thus,  $fop_i = g^{(\sum_{j=1}^{j=m} (x_i)_j - m)} \bmod \eta$ . Only when  $\sum_{j=1}^{j=m} (x_i)_j = m$ , the result of above expression is one. Otherwise, it is a nonzero number.

**Information leakage discussion.** We need to prevent information leakage at the server and at the DB owners.

**1) Server perspective.** The servers only know the parameters  $\langle g, \delta, \eta' \rangle$  and may utilize the relations between  $g$  and  $\eta$  to guess  $\eta$  from  $\eta'$ . However, it will not give any meaningful information to servers, since the DB owner sends the elements of  $\chi$  in additive shared form, and since servers do not know each other, they cannot obtain the cleartext values of  $\chi$ . Also, an identical operation is executed on all shares of  $m$  DB owners. Hence, access-patterns are hidden from servers, thereby the server cannot distinguish between any two values based on access-patterns. Further, the output of query is also in shared form and contains an identical number of bits. Thus, based on the output size, the server cannot know whether the value is common among DB owners or not.

**2) DB owner perspective.** When all DB owners do not have one at the  $i^{th}$  position of  $\chi$ , we need to inform DB owners that there is no common value and not to reveal that how many DB owners do not have one at the  $i^{th}$  position. Note that the DB owner can learn this information, if they know  $g$  and  $\alpha$ , since based on these values, they can compute what the servers have computed. However, unawareness of  $g$  and  $\alpha$  makes it impossible to guess the number of DB owners that do not have one at the  $i^{th}$  position of  $\chi$ . We can formally prove it as follows:

**Lemma.** A DB owner cannot deduce how many other DB owners do not have one at the  $i^{th}$  position of  $\chi$  without knowing  $g$ .

**Proof.** According to the precondition,  $g$  is a generator of a cyclic group of order  $\delta$ , where  $\delta$  is a prime number. Thus,  $C = \{g^0, g, g^2, \dots, g^{\delta-1}\}$  represents all items in the cyclic group. Assume that the output of Equation 4 is a number other than one, say  $\beta$ . Thus, we have  $\beta = g^{x-m} \bmod \eta$ , where  $x$  represents the number of one at the  $i^{th}$  position of  $\chi_j$ ,  $1 \leq j \leq m$ . When DB owners wish to know  $x$ , they must compute  $\log_g \beta$ . To solve it, they need to know  $g$ . Note that based on the characteristic of the cyclic group, there are less than  $\delta-1$  generators of  $C$  and co-prime to  $\delta$ . Thus,  $g^2, \dots, g^{\delta-1}$  may also be generators of the cyclic group. However, DB owners cannot distinguish which generator is used by the servers. Thus, DB owners cannot deduce the value of  $x$ , except knowing that  $x \in [0, m-1]$ .  $\blacksquare$

## 5.2 PSI Result Verification

A malicious adversary or a hardware/software bug may result in the following situations, during computing PSI: (i) skip processing the  $i^{th}$  additive shares of all DB owners, (ii) replacing the result of the  $i^{th}$  additive shares by the computed result for  $j^{th}$  share, (iii) injecting fake values, or (iv) falsifying the verification method. Thus, this section provides a method for verifying the result of PSI.

**High-level idea.** Let  $g$  be a generator of a cyclic group under modulo multiplicative  $\eta$  operation, and  $\eta' = \alpha \times \eta$ ,  $\alpha > 1$ . Thus,  $(g^x \bmod \eta) \times (g^{-x} \bmod \eta) = 1$ , and the idea of PSI verification lies in this equation. Recall that in PSI (§5.1), we used  $(g^x \bmod \eta)$ , whose value 1 shows that the item exists at all DB owners. Now, we will use the term  $(g^{-x} \bmod \eta)$  for verification. Specifically, if the servers have performed their computations correctly, then Equation 5 must hold to be true:

$$((g^{(\oplus_{j=1}^{j=m} A(x_i)_j^1) - A(m)^1} \bmod \eta') \times (g^{(\oplus_{j=1}^{j=m} A(x_i)_j^2) - A(m)^2} \bmod \eta')) \bmod \eta = 1 \quad (5)$$

where  $m$  is the number of DB owners,  $x_j$  is either 1 or 0 (as described in §5.1), and  $\bar{x}_j$  is the complement value of  $x_j$ . Below, we describe the steps executed at the servers and DB owners.

**STEP 1: DB owners.** On the distinct values of an attribute  $A_c$  of their relations, the DB owner  $\mathcal{DB}_j$  executes a hash function to create the table  $\chi_j$  that contains  $b = |\text{Dom}(A_c)|$  values (either 0 or 1). Further,  $\mathcal{DB}_j$  creates a table  $\bar{\chi}_j$  containing  $b$  values, such that  $i^{th}$  value of  $\bar{\chi}_j$  must be the complement of  $i^{th}$  value of  $\chi_j$ . Then,  $\mathcal{DB}_j$  permutes the values of  $\bar{\chi}_j$  using a permutation function  $\mathcal{PF}_{db1}$  (known to all DB owners only) and creates additive shares of each value of  $\chi_j$  and  $\bar{\chi}_j$ , prior to outsourcing to servers. The reason of using  $\mathcal{PF}_{db1}$  will be clear soon.

**STEP 2: Servers.** Each server  $\mathcal{S}_\phi$  holds the  $\phi^{th}$  additive share of  $\chi$  (denoted by  $A(\chi)_j^\phi$ ) and  $\bar{\chi}$  (denoted by  $A(\bar{\chi})_j^\phi$ ) of  $j^{th}$  DB owner and executes the following operation:

$$output_i^{S_\phi} \leftarrow g^{((\oplus_{j=1}^{j=m} A(x_i)_j^1) \ominus A(m)^1) \bmod \eta', (1 \leq i \leq b)} \quad (6)$$

$$Vout_i^{S_\phi} \leftarrow g^{((\oplus_{j=1}^{j=m} A(\bar{x}_i)_j^1) \ominus A(m)^1) \bmod \eta', (1 \leq i \leq b)} \quad (7)$$

Equation 6 is used to find the common item at the server and is identical to Equation 3, given in §5.1. In Equation 7, each server  $\mathcal{S}_\phi$  performs the following operations: (i) modular addition (under  $\delta$ ) of the  $i^{th}$  additive shares of  $\bar{\chi}$  from  $m$  DB owners, (ii) exponentiation by  $g$  to the

<sup>1</sup> Consider  $i^{th}$ ,  $j^{th}$ , and  $k^{th}$  values of  $\chi_1 = \{1, 0, 0\}$ ,  $\chi_2 = \{0, 1, 0\}$ ,  $\chi_3 = \{1, 1, 1\}$ . Here, after STEP 3, DB owners will learn three random numbers, such that the first two random numbers will be identical. Based on this, DB owner can only know that the sum of  $i^{th}$  and  $j^{th}$  position of  $\chi$  is identical. However, it will not reveal how many positions have 0 or 1 at  $i^{th}$  or  $j^{th}$  positions.



Value	Share 1	Share 2	Value	Share 1	Share 2	Value	Share 1	Share 2
0	2	-2	0	2	-2	0	4	-4
1	0	1	0	3	-3	1	1	0
0	1	-1	1	4	-3	0	1	-1

Table 9:  $\mathcal{DB}_1$ .Table 10:  $\mathcal{DB}_2$ .Table 11:  $\mathcal{DB}_3$ .

power the result of the previous step, under modulus  $\eta'$ ; and (iii) sends the computed results  $output^{S_\phi}[]$  and  $Vout^{S_\phi}[]$  to all DB owners.

**STEP 3: DB owners.** From two servers, DB owners receive  $output^{S_\phi}[]$  and  $Vout^{S_\phi}[]$  (each of length  $b$ ), permute back the values of  $Vout^{S_\phi}[]$  (using the reverse permutation function, since they used  $\mathcal{PF}_{db1}$  on  $\bar{\chi}$ , which results in  $Vout^{S_\phi}[]$  at servers) to obtain  $pvout^{S_\phi}[]$ , and execute the following:

$$r_1 \leftarrow output_i^{S_1} \times output_i^{S_2} \bmod \eta \quad (8)$$

$$r_2 \leftarrow pvout_i^{S_1} \times pvout_i^{S_2} \bmod \eta \quad (9)$$

$$r_1 \times r_2 \bmod \eta = 1 \quad (10)$$

If the DB owner find the output of  $r_1 \times r_2$  equals to one for all  $b$  values, it shows that the servers executed the computation correctly.

**Example 5.2.1.** We verify PSI results of Example 5.1.1. Suppose  $\delta=5$ ,  $\eta=11$ , and  $\eta'=143$ , as assumed in Example 5.1.1.

**Step 1: DB owners.** DB owners find the reverse of  $\chi$  (as shown in the first column of Tables 9, 10, and 11) and generate additive shares; see the second and third columns of Tables 9, 10, and 11. Note that here for simplicity, we do not permute the values or shares.

**Step 2: Servers.** The server  $S_1$  will return the three values 27, 81, 3, by executing the following computation, to all three DB owners:

$$\begin{aligned} 3^{((2+2+4) \bmod 5)} \bmod 143 &= 27 \\ 3^{((0+3+1) \bmod 5)} \bmod 143 &= 81 \\ 3^{((1+4+1) \bmod 5)} \bmod 143 &= 3 \end{aligned}$$

The server  $S_2$  will return three values 7, 27, and 1 to all three DB owners:

$$\begin{aligned} 3^{((-2-2-4) \bmod 5)} \bmod 143 &= 9 \\ 3^{((1-3+0) \bmod 5)} \bmod 143 &= 27 \\ 3^{((-1-3-1) \bmod 5)} \bmod 143 &= 1 \end{aligned}$$

**Step 3: DB owners.** The DB owner obtains a vector  $\langle 1, 9, 8 \rangle$ , by executing the following computation:

$$\begin{aligned} (27 \times 9) \bmod 11 &= 1 \\ (81 \times 27) \bmod 11 &= 9 \\ (3 \times 1) \bmod 11 &= 3 \end{aligned}$$

Now, the DB owner executes the following for verifying PSI results,  $1 \times 1 \bmod 11 = 1$ ,  $5 \times 9 \bmod 11 = 1$ , and  $4 \times 3 \bmod 11 = 1$ , where 1, 5, 4 are the final output at DB owner in Example 5.1.1. The output 1 indicates that the servers executed the computation correctly. ■

**Correctness.** First, we need to argue that the processing at servers works correctly. Assume that the DB owner does not implement  $\mathcal{PF}_{db1}$  on elements of  $\bar{\chi}$ , and computation at servers is executed in cleartext. Thus, on the values of  $\chi$ , servers add  $i^{th}$  value of each  $\chi_j = \{x_i\}$  ( $1 \leq j \leq m$ ,  $1 \leq i \leq b$ ) and subtract the results from  $m$ . It will result in a number, say  $a \in \{-m+1, 0\}$ . On the other hand, servers add  $i^{th}$  values  $\bar{\chi}_j$ , and it will result in a number, say  $b \in \{0, m\}$ , i.e., the number of ones at DB owners at the  $i^{th}$  position of  $\chi$ . To hide the value of  $a$  and  $b$  from servers, they execute operations on additive shares of  $\chi$  and  $\bar{\chi}$ , and take a modulus exponent (i.e.,  $r_1 \leftarrow g^a$  and  $r_2 \leftarrow g^b$ ) to hide  $a$  and  $b$  from DB owners. Since  $a = -b$  or  $a = b = 0$ ,  $r_1 \times r_2 \bmod \eta = 1$ , and this shows that the server executed the correct operation.

Now, we show why the verification method will detect any abnormal computation executed by servers. Note that servers may skip to process all/some values of  $\chi$  and  $\bar{\chi}$ . For example, servers may process only  $x_1 \in \chi$ ,  $\bar{x}_1 \in \bar{\chi}$ , and send the results corresponding to  $x_1$ ,  $\bar{x}_1$  as the results of all

remaining  $b-1$  values. Such a malicious operation of servers will provide a legal proof (i.e.,  $r_1 \times r_2 \bmod \eta = 1$ ) at DB owners that servers executed the computation correctly, (since values of  $\bar{\chi}$  was not permuted). Thus, we used permutation over the values of  $\bar{\chi}$  and/or additive shares of  $\bar{\chi}$ . Now, to break the verification method and to produce  $r_1 \times r_2 \bmod \eta = 1$  for an  $i^{th}$  value of  $\chi$ , servers need to find the correct value in  $\bar{\chi}$  corresponding to an  $i^{th}$  value of  $\chi$  (among the randomly permuted shares). Hence, the removal of any results from the output will be detected.

Now, we show why the verification method can detect fake data insertion by servers. For a malicious server  $S_1$  to successfully inject a fake tuple (i.e., undetected during verification), it should know the correct position of some element in both  $A(\chi)_j^1$  and  $A(\bar{\chi})_j^1$ . Since  $A(\bar{\chi})_j^1$  is a permuted vector of size  $b = |\text{Dom}(A_c)|$ , the probability of finding the correct element in  $A(\bar{\chi})_j^1$  corresponding to an element of  $A(\chi)_j^1$  will be  $1/b^2$ . E.g., in our experiments, the domain size is 5M (or 20M) values, making the above probability infinitesimal ( $< 10^{-13}$ ).<sup>2</sup> **Additional security.** We implemented  $\mathcal{PF}_{db1}$  on the elements of  $\bar{\chi}$ . We can, further, permute additive shares of both  $\chi$  and  $\bar{\chi}$  using different permutation functions, to make it impossible for both servers to find the position of a value in  $A(\chi)_j^\phi$  and  $A(\bar{\chi})_j^\phi$ ,  $\phi \in \{1, 2\}$ . Thus, servers cannot break the verification method, and any malicious activities will be detected by DB owners.

**Information leakages discussion.** Arguments for information leakage follows the similar way as the arguments given for PSI computation in §5.1. Thus, the verification method will not reveal any non-desired information to servers and DB owners.

## 6 AGGREGATION OPERATION OVER PSI

PRISM supports both summary and exemplar aggregations. Below, we describe how PRISM implements sum §6.1, average §6.2, maximum §6.3, median §6.4 and count operations §6.5. Also, in our discussion below, we will consider set-based operation PSI on a single attribute  $A_c$ . §6.6 will extend the discussions to support PSI over on multiple attributes and over a large-size domain.

### 6.1 PSI Sum Query

A PSI sum query computes the sum of values over an attribute corresponding to common items in another attribute; see example given in §2. This section develops a method based on additive, as well as, multiplicative shares, where additive shares find common items over an attribute  $A_c$  and multiplicative shares (SSS) finds the sum of shares of an attribute  $A_x$  corresponding to the common items in  $A_c$ . This method contains the following steps:

**STEP 1: DB owners.**  $\mathcal{DB}_j$  creates their  $\chi_j$  table over the distinct values of  $A_c$  attribute by following STEP 1 of PSI; see §5. Here,  $\chi_j = \{\langle x_{i1}, x_{i2} \rangle\}$ , where  $1 \leq i \leq b$  and  $b = |\text{Dom}(A_c)|$ , i.e., the  $i^{th}$  cell of  $\chi_j$  contains a pair of values,  $\langle x_{i1}, x_{i2} \rangle$ , where (i)  $x_{i1} = 1$ , if a value  $a_i \in A_c$  is mapped to the  $i^{th}$  cell, otherwise, 0; and (ii)  $x_{i2}$  contains the sum of values of  $A_x$  attribute corresponding to  $a_i$ ; otherwise, 0.  $\mathcal{DB}_j$  creates additive shares of  $x_{i1}$  (denoted by  $A(x_{i1})_j^\phi$ ,  $\phi = \{1, 2\}$ ) and sends to two servers  $S_1$  and  $S_2$ . Also,  $\mathcal{DB}_j$  creates SSS of  $x_{i2}$  (denoted by  $S(x_{i2})_j^\phi$ ,  $\phi = \{1, 2, 3\}$ ) and sends to three servers  $S_1$ ,  $S_2$ , and  $S_3$ .

**STEP 2: Servers.** Servers  $S_1$  and  $S_2$  find common items using additive shares by implementing Equation 3 and send all computed

<sup>2</sup>If the domain size is small, we can increase its size by adding fake values to bind the probability of adversary being able to inject fake data.

$b$  results to all DB owners. Since the result is in additive shared form, it cannot be multiplied to SSS. Thus, **servers send the output of PSI to one of the DB owners selected randomly and wait to receive multiplicative shares corresponding to common items. The reason of randomly selecting only one DB owner is just to reduce the communication overhead of sending/receiving additive/multiplicative shares, and it does not impact the security. Note that all DB owners can receive the PSI outputs and generate multiplicative shares.**

**STEP 3: DB owners.** On receiving  $b$  values, the DB owner finds the common items by executing Equation 4 and generates a vector of length  $b$  having 1 or 0 only, where 0 is obtained by replacing random values of  $fop$ . Finally, the DB owner creates three SSS of each of the  $b$  value, denoted by  $S(z_i)^\phi$ ,  $\phi = \{1,2,3\}$ , and sends to three servers.

**STEP 4: Servers.** Servers  $S_\phi$ ,  $\phi = \{1,2,3\}$ , execute the following:

$$sum_i^{S_\phi} \leftarrow \sum_{j=1}^{j=m} S(x_{i2})_j^\phi \times S(z_i)^\phi, 1 \leq i \leq b \quad (11)$$

Each server multiplies  $S(z_i)^\phi$  by  $S(x_{i2})_j^\phi$  of each DB owner, adds the results, and sends them to all DB owners.

**STEP 5: DB owners.** From three servers, all DB owners receive three vectors, each of length  $b$ , and perform Lagrange interpolation on each  $i^{th}$  value of the three vectors to obtain the final sum of the value in  $A_x$  corresponding to the common items in  $A_c$  attribute, which they received in STEP 3.

**Correctness.** Here, we need to show the correctness of PSI, which can be argued similarly as presented in §5. PSI results in one (if an item is common at all DB owners) or random numbers (otherwise) that will be replaced by zero. Such 0 or 1 values are converted into multiplicative shares using SSS in STEP 3. Since servers multiply an  $i^{th}$  SSS (which is zero or one) in STEP 4 to the  $i^{th}$  SSS of  $A_x$  and then add all  $i^{th}$  values of  $A_x$ , it will result in either zero or the sum of shares sent by DB owner. Thus, servers can compute the sum of shares corresponding to common items among DB owners.

**Information leakage.** By following the argument of information leakage of PSI (see §5.1), we can state here that servers cannot learn the intersection and values. Now, we can argue about information leakage due to the sum of values, as follows: since the values of  $A_x$  attribute are in SSS form, servers cannot learn the actual value. During multiplication operations, servers perform identical operations on each SSS; thus, based on access-patterns for finding sum corresponding to common items, the adversary neither learns the common item nor the sum of shares. Since DB owners only receive the result of multiplication of SSS, DB owners cannot learn anything more than the sum of values of  $A_x$  corresponding to the common values in  $A_c$ .

**PSI sum verification.** PSI sum verification requires to verify two things: the common items in  $A_c$  and the sum of value of  $A_x$  corresponding to the common items. We verify common items using PSI verification method §5.2. For sum verification, we do the following:

**STEP 1: DB owner.**  $\mathcal{DB}_j$  creates three vectors: (i)  $\chi_j = \{\langle x_{i1}, x_{i2} \rangle\}$  as in STEP 1 of PSI Sum, (ii)  $\overline{\chi_j} = \{\langle \overline{x_{i1}} \rangle\}$ , i.e.,  $\overline{\chi_j}$  contains complement of each value  $x_{i1}$  in a permuted order using a permutation function  $\mathcal{PF}_{db1}$ , as we did in STEP 1 of PSI verification, and (iii)  $\chi'_j = \{\langle px_{i2} \rangle\}$  that contain  $x_{i2}$  of  $\chi_j$  in a permuted order using a permutation function  $\mathcal{PF}_{db2}$ . Recall that  $\mathcal{PF}_{db1}$  and  $\mathcal{PF}_{db2}$  are only known to all DB owners. Here,  $x_{i1}$  and  $\overline{x_{i1}}$  are used to verify PSI, as explained in §5.2. The values  $x_{i2}$  and  $\chi'_j$  are used to verify the sum.  $\mathcal{DB}_j$  creates SSS of  $\langle x_{i2}, px_{i2} \rangle$  and additive shares of  $\langle x_{i1}, \overline{x_{i1}} \rangle$ .

**STEP 2: Servers.** Servers find PSI, execute PSI verification method by Equations 3,6,7, and send two output vectors each of length  $b$  to all DB owners.

**STEP 3: DB owners.** DB owners verify PSI output using Equations 8-10. Additionally, **only one of the DB owners** generates two vectors  $\Gamma_1$  and  $\Gamma_2$ , each of length  $b$  having 1 or 0 only, corresponding to common items outputs. Also, it permutes  $\Gamma_2$  using  $\mathcal{PF}_{db2}$ , and creates SSS of each value in both vectors, prior to send to three servers. We denote SSS of  $\gamma_i \in \Gamma_1$  by  $S(\gamma_i)^\phi$  and SSS of  $\rho \in \Gamma_2$  by  $S(\rho_i)^\phi$ .

**STEP 4: Servers.** Servers  $S_\phi$ ,  $\phi = \{1,2,3\}$ , execute the following:

$$sum_i^{S_\phi} \leftarrow \sum_{j=1}^{j=m} S(x_{i2})_j^\phi \times S(\gamma_i)^\phi, 1 \leq i \leq b \quad (12)$$

$$vsum_i^{S_\phi} \leftarrow \sum_{j=1}^{j=m} S(px_{i2})_j^\phi \times S(\rho_i)^\phi, 1 \leq i \leq b \quad (13)$$

**STEP 5: DB owners.** DB owners perform Lagrange interpolation to obtain the final sum of the value in  $A_x$ . Also, they permute back the values of the permuted vector (i.e., all values of  $vsum_i^{S_\phi}$  after interpolation) and match against the output of the non-permuted vector (i.e., output of  $vsum_i^{S_\phi}$  against  $sum_i^{S_\phi}$  after interpolation). If both vectors' outputs match, it shows servers have executed the computation correctly.

**Correctness.** As DB owners outsource multiplicative shares of  $x_{i2}$  with  $px_{i2}$  ( $1 \leq i \leq b$ ), where the shares  $px_{i2}$  in the vector are permuted, to void the sum verification either by executing wrong computation or by injecting fake data, servers need to know the correct position of some elements in both vectors having multiplicative shares. Otherwise, a wrong computation or fake tuples insertion at any position, say  $i$ , will result in non-identical values of  $sum_i^{S_\phi}$  and  $vsum_i^{S_\phi}$ , after interpolation and permutation of the values. Thus, sum verification method will detect fake tuple insertion or wrong computation by servers.

## 6.2 PSI Average Query

A PSI average query on cost column corresponding to the common disease in Tables 1-3 returns {Cancer, 280}. PSI average query works in a similar way as PSI sum query. In short,  $\mathcal{DB}_j$  creates  $\chi_j = \{\langle x_{i1}, x_{i2}, x_{i3} \rangle\}$ , where  $1 \leq i \leq b$ ,  $b = |\text{Dom}(A_c)|$ , and  $x_{i1}, x_{i2}$  are identical to the values we created in STEP 1 of PSI sum (§6.1). The new value  $x_{i3}$  contains the number of tuples at  $\mathcal{DB}_j$  corresponding to  $x_{i1}$ . For example, in case of Table 1, one of the values of  $\chi_1$  will be  $\{\langle \text{Cancer}, 300, 2 \rangle\}$ , i.e., Table 1 has two tuples corresponding to Cancer and cost 300. All values  $x_{i3}$  are outsourced in multiplicative share form. Then, we follow STEPS 2 and 3 of PSI sum. In STEP 4, the servers also multiply the received  $i^{th}$  SSS values corresponding to the common value to  $x_{i2}, x_{i3}$  and add the values. Finally, in STEP 5, DB owners interpolate vectors corresponding to all  $b$  values of  $x_{i2}, x_{i3}$  and find the average by dividing the values appropriately.

**Correctness:** can be argued similar to PSI sum.

**Information leakage:** can be argued similar to PSI sum. Note that here we reveal the total number of tuples and the sum of values corresponding to the common values.

## 6.3 PSI Maximum Query

This section develops a method for finding the maximum value in an attribute  $A_x$  corresponding to the common values in  $A_c$  attribute; refer to §2 for PSI maximum example. Here, our objective is to prevent the adversarial server from learning: (i) the actual maximum values outsourced by each DB owner, (ii) what is the maximum value among DB owners and which DB owners have the maximum value. We allow all the DB owners to know the maximum value and/or the identity



of the DB owner(s) having the maximum value. We use pick color to highlight the part that is used to reveal the identity of DB owners having maximum to distinguish which part of the algorithm can be avoided based on the security requirements.

In this method, *each DB owner uses the polynomial  $\mathcal{F}(x)$*  given by the initiator; see §4 to find how we created  $\mathcal{F}(x)$ . Note that we use this polynomial to generate values at different DB owners in an order-preserving manner by executing the following STEP 3 and Equation 14.

The method contains at most three rounds, where the first round finds the common values in an attribute  $A_c$  by using STEPS 1-3, the second round finds the maximum value in an attribute  $A_x$  corresponding to common items in  $A_c$  using STEPS 4-5a, the last round finds DB owners who have the maximum value using STEPS 5b-7. Note that *the third round is not always required*, if (i) we do not want to reveal the identity of the DB owner having the maximum value, or (ii) values in  $A_x$  column across all DB owners are unique.

**STEP 1 at DB owner and STEP 2 at servers.** These two steps are identical to STEP 1 and STEP 2 of PSI query execution method (§5).

**STEP 3: DB owner.** On the received outputs (of STEP 2) from servers, DB owners find the common item in the attribute  $A_c$ , as in STEP 3 of PSI query execution method (§5). Now, to find the maximum value in the attribute  $A_x$  corresponding to the common item in  $A_c$ , DB owners proceeds as follows:

For the purpose of simplicity, we assume that there is only one common item, say  $y^{th}$  item.  $\mathcal{DB}_i$  finds the maximum, say  $M_{iy}$ , in the attribute  $A_x$  of their relation corresponding to the common item  $y$ . Note that since we assume only one common element, we refer to the maximum element  $M_{iy}$  by  $M_i$ .  $\mathcal{DB}_i$  executes Equation 14 to produce values at DB owners in an order-preserving manner:

$$v_i \leftarrow \mathcal{F}(M_i) + r_i \quad (14)$$

$\mathcal{DB}_i$  implements the polynomial  $\mathcal{F}()$  on  $M_i$  and adds a random number  $r_i$  (selected in a range between 0 and  $M_i^m$ ), and it produces a value  $v_i$ . Finally,  $\mathcal{DB}_i$  creates additive shares of  $v_i$  (denoted by  $A(v_i)^\phi$ ) and sends them to servers  $\mathcal{S}_\phi$ ,  $\phi \in \{1, 2\}$ . Note that even if  $k \geq 2$  DB owners have the same maximum value  $M_i$ , by this step, the value  $v$  will be different at those DB owners, with a high probability,  $1 - \frac{1}{(M_i)^{(k-1)m}}$ , (depending on the range of  $r_i$ ). Also, if any two numbers  $M_i < M_j$ , then  $\mathcal{F}(M_i) + r_i < \mathcal{F}(M_j)$  will hold. Full version will prove such statements.

**STEP 4: Servers.** Each server  $\mathcal{S}_\phi$  executes the following operation:

$$input^{\mathcal{S}_\phi}[i] \leftarrow A(v_i)^\phi, 1 \leq i \leq m; output^{\mathcal{S}_\phi}[] \leftarrow \mathcal{PF}(input^{\mathcal{S}_\phi}[])$$

Server  $\mathcal{S}_\phi$  collects additive shares from each DB owner and places them in an array (denoted by  $input^{\mathcal{S}_\phi}[]$ ), on which  $\mathcal{S}_\phi$  executes the permutation function  $\mathcal{PF}$ . Then, they send the output the permutation function, i.e.,  $output^{\mathcal{S}_\phi}[]$ , to the announcer  $\mathcal{S}_a$  that executes the following:

$$foutput^{\mathcal{S}_a}[i] \leftarrow output^{\mathcal{S}_1}[i] + output^{\mathcal{S}_2}[i], 1 \leq i \leq m \quad (15)$$

$$max, index \leftarrow FindMax(foutput^{\mathcal{S}_a}[]) \quad (16)$$

$\mathcal{S}_a$  adds the  $i^{th}$  outputs received from  $\mathcal{S}_1$  and  $\mathcal{S}_2$ , and compares all those numbers to find the maximum number (denoted by  $max$ ). Also,  $\mathcal{S}_a$  produces the index position (denoted by  $index$ ) corresponding to the maximum number in  $foutput^{\mathcal{S}_3}[]$ . Finally,  $\mathcal{S}_a$  creates additive secret-shares of  $max$  (denoted by  $A(max)^{\mathcal{S}_\phi}$ ,  $\phi \in \{1, 2\}$ ), as well as, of  $index$  (denoted by  $A(index)^{\mathcal{S}_\phi}$ ), and sends them to  $\mathcal{S}_\phi$  ( $\phi \in \{1, 2\}$ ) that forwards such additive shares to DB owners. Note that *if the protocol*

*does not require to reveal the identity of the DB owner having the maximum value,  $\mathcal{S}_a$  does not send additive shares of  $index$ .*

**STEP 5a: DB owner.** Now, the DB owners' task is to find the maximum value and/or the identity of the DB owner who has the maximum value. To do so, each DB owner performs the following:

$$max \leftarrow A(max)^{\mathcal{S}_1} + A(max)^{\mathcal{S}_2} \quad (17)$$

$$index \leftarrow A(index)^{\mathcal{S}_1} + A(index)^{\mathcal{S}_2}, pos \leftarrow \mathcal{RPF}(index) \quad (18)$$

The DB owner finds the identity of the DB owner having the maximum value by adding the additive shares and by implementing reverse permutation function  $\mathcal{RPF}$ . Note that  $\mathcal{RPF}$  works since  $\mathcal{PF}$  is known to DB owners and servers (see Assumptions given in §4). To find the maximum value, they implement  $\mathcal{F}(z)$  and  $\mathcal{F}(z+1)$  and evaluate  $\mathcal{F}(z) \leq max < \mathcal{F}(z+1)$ , where  $z \in \{1, 2, \dots\}$ .<sup>3</sup> If this condition holds to be true, then  $z$  is the maximum value, and if  $z = M_i$ , then the  $i^{th}$  DB owner knows that he/she holds the maximum value. Obviously, if the  $i^{th}$  DB owner does not hold the maximum value, then  $M_i < \mathcal{F}(M_i) + r_i < \mathcal{F}(M_i + 1) \leq \mathcal{F}(z) \leq max$ .

**STEP 5b: DB owner.** Note that by the end of STEP 5a, the DB owners know the maximum value and the identity of the DB owner having the same maximum value, due to  $pos$ . However, if there are more than one DB owner having the maximum value, the other DB owners cannot learn about it. The reason is: the server  $\mathcal{S}_a$  can find only the maximum value, while, recall that, if more than one DB owners have the same maximum value, say  $M$ , they produce a different value, due to using different random numbers in STEP 3 (Equation 14). Thus, we need to execute this step 5b to know all DB owners having the maximum value.

After comparing its maximum values against  $max$ ,  $\mathcal{DB}_i$  knows whether it possesses the maximum value or not. Depending on this,  $\mathcal{DB}_i$  generates a value  $\alpha_i = 0$  or  $\alpha_i = 1$ , creates additive shares of  $\alpha_i$ , and sends to  $\mathcal{S}_\phi$ ,  $\phi \in \{1, 2\}$ .

**STEP 6: Servers.** Server  $\mathcal{S}_\phi$  allocates the received additive shares to a vector, denoted by  $fpos$ , and sends the vector  $fpos$  to all DB owners, i.e.,  $fpos^{\mathcal{S}_\phi}[i] \leftarrow A(\alpha_i)^{\mathcal{S}_\phi}$ ,  $1 \leq i \leq m$ .

**STEP 7: DB owner.** Each DB owner adds the received additive shares to obtain the vector  $fpos[]$ .

$$fpos[i] \leftarrow fpos^{\mathcal{S}_1}[i] + fpos^{\mathcal{S}_2}[i], 1 \leq i \leq m \quad (19)$$

By  $fpos[]$ , DB owners discover which DB owners have the maximum value, since, recall that in STEP 5,  $\mathcal{DB}_i$  that satisfies the condition  $(\mathcal{F}(M_i) \leq max < \mathcal{F}(M_i + 1))$  requests  $\mathcal{S}_\phi$  to place additive share of 1 at  $fpos^{\mathcal{S}_\phi}[i]$ .

**Example 6.3.1.** Assume  $\eta = 5003$ . Refer to Tables 1-3, and consider that all hospitals wish to find the maximum age of a patient corresponding to the common disease and which hospitals treat such patients. We assume that all hospitals know cancer as the common disease.

In STEP 3, all hospitals, i.e., DB owners, find their maximum values in the attribute Age corresponding to common disease and implement  $\mathcal{F}(x) = x^4 + x^3 + x^2 + x + 1$ , sent by the initiator.

$$\mathcal{F}(6) = 1555 + 216 = 1771 = (5000 - 3229) \bmod 5003$$

$$\mathcal{F}(8) = 4681 + 1 = 4682 = (5500 - 818) \bmod 5003$$

$$\mathcal{F}(8) = 4681 + 319 = 5000 = (2500 + 2500) \bmod 5003$$

Further, they add random numbers (216, 1, 319) and create additive shares, which are outsourced to  $\mathcal{S}_1$  and  $\mathcal{S}_2$ . In STEP 4,  $\mathcal{S}_1$  holds  $\langle 5000, 5500, 2500 \rangle$ , permutes them, and sends to  $\mathcal{S}_a$ .  $\mathcal{S}_2$  holds  $\langle -3229, -818, 2500 \rangle$ , permutes them, and sends to  $\mathcal{S}_a$ .

<sup>3</sup>To reduce the computation cost, we can select number  $z$  similar to binary search method.

$S_a$  obtains  $\langle 4682, 5000, 1771 \rangle$  by adding the received shares from  $S_1, S_2$ , and finds 5000 as the maximum value and ‘Hospital 2’ to which the value belongs. Finally,  $S_a$  creates additive shares of  $5000 = (4000 + 1000) \bmod 5003$ , **additive shares of the identity of the DB owner as  $2 = (200 - 198) \bmod 5003$** , and sends to DB owners via  $S_1, S_2$ .

In STEP 5a, all hospitals will know the maximum value as 5000 (with random value added) **and identity of the DB owner as 2 on which they implement the reverse permutation function to obtain the correct identity as ‘Hospital 3’**. Then, ‘Hospital 1’ knows that they do not hold the maximum, since  $\mathcal{F}(6) + 216 < \mathcal{F}(7) < 5000$ . ‘Hospital 2’ knows that they hold the maximum, since  $\mathcal{F}(8) < 5000 < \mathcal{F}(9)$ . Also, ‘Hospital 3’ knows that they hold the maximum.

**To know which hospitals have the maximum value, in STEP 5b, Hospitals 1, 2, 3’ create additive shares of 0, 1, 1, respectively, as:  $0 = (200 - 200) \bmod 5003$ ,  $1 = (300 - 299) \bmod 5003$ , and  $1 = (200 - 199) \bmod 5003$ , and send to  $S_1$  and  $S_2$ . Finally, in STEP 6,  $S_1$  and  $S_2$  send  $\langle 200, 300, 200 \rangle$  and  $\langle -200, -299, -199 \rangle$  to all hospitals. In STEP 7, hospitals add received shares, resulting in  $(0, 1, 1)$ . It shows that ‘Hospitals 2, 3’ have the maximum value 8. ■**

**Correctness.** We need to show the proposed method will allow: (i) the DB owners to know the maximum value, and (ii) **the identity of the DB owners who have the maximum value.**

First, before showing that the DB owners will know the maximum value, we need to show that the server will find the maximum value. Recall that in STEP 3, the random number addition to output of  $\mathcal{F}(\cdot)$  hides only the actual value from  $S_a$ , (i.e., if  $a < b$ , then  $\mathcal{F}(a) + r < \mathcal{F}(b) + r'$  and if  $a = b$ , then  $\mathcal{F}(a) + r < \mathcal{F}(b) + r'$  or  $\mathcal{F}(a) + r > \mathcal{F}(b) + r'$  will hold at  $S_a$ , where  $r, r'$  are random numbers). Also, observe that the execution of the permutation function by  $S_1$  and  $S_2$  just hides the identity of DB owners from  $S_a$  and does not affect the shares. Thus, in STEP 5,  $S_a$  will find the maximum value (among the permuted shares, received from  $S_1$  and  $S_2$ ) and provide additive shares of the maximum value and **additive shares of the identity of the DB owner.**

Second, we show that each DB owner will know the maximum value — by comparing her maximum value from the summation of the maximum value’s additive shares, received from  $S_1$  and  $S_2$ . Consider three values  $a < b = c$  at three DB owners  $\mathcal{DB}_1, \mathcal{DB}_2, \mathcal{DB}_3$ , respectively, and these values become  $v_1 < v_2 < v_3$  using STEP 3, Equation 14. Here, of course,  $S_a$  will announce  $v_3$  as the maximum value, and thus, on receiving additive shares of the maximum value, via  $S_1$  and  $S_2$ ,  $\mathcal{DB}_3$  will know it has the maximum values. Also,  $\mathcal{DB}_2$  will know that it has the maximum value, since  $\mathcal{F}(b) \leq v_3 < \mathcal{F}(b+1)$  (see STEP 5a). Also,  $\mathcal{DB}_1$  will know that it does not hold the maximum value, since  $\mathcal{F}(a) < \mathcal{F}(a+1) \leq v_3$ .

**Finally, once all DB owners will know whether they hold the maximum value or not, using STEPS 5a-7, they will also know the identity of each DB owner who has the maximum value, by checking the value of  $fpos[]$  (STEP 7), since  $\mathcal{DB}_i$  with the maximum value has requested  $S_\phi$  to place additive shares of 1 at the  $fpos^{S_\phi}[i]$  (STEP 5b).**

**Information leakage discussion.** We discuss information leakage at the servers and at the DB owner.

• **Server perspective.** Here, our objective is to hide: (i) the actual values, (ii) the number of DB owners having the same value, (iii) the maximum value and the identity of the DB owners who have the maximum value from the servers, and (iv)  $S_a$  cannot reconstruct the actual value, since it receives the output of STEP 3. In order to achieve the first two objectives,

in STEP 3, the DB owners use a polynomial and add a random number to the output of the polynomial. Thus, even if two or more DB owners have the same or different value, their final output value will be different.

Since  $S_a$  finds the maximum value and the identity over the permuted values,  $S_a$  cannot deduce that which value is related to which DB owner. Since  $S_a$  sends additive shares of the maximum value and identity to  $S_1$  and  $S_2$ , they cannot learn that which DB owner has the maximum value. Thus, we achieve the third objective.

In order to achieve the fourth objective, the DB owners use a polynomial of degree more than  $m$ , where  $m$  is the number of DB owners. Thus, even collecting values from all  $m$  DB owners,  $S_a$  cannot interpolate the received values to know the actual value. Observe that this statement is akin to Shamir’s secret-sharing, where an adversary cannot learn a secret, until collecting  $t+1$  shares, if a polynomial of degree  $t$  is used for creating shares of a secret.

• **DB owner perspective.** DB owners’ objectives are (i) the servers will not learn their actual values  $\mathcal{M}$ , (ii) each DB owner will not learn other DB owners’ values, except the maximum value, and (iii) if they are not interested in learning the identity of the DB owner having the maximum value, it should not be revealed. Due to not revealing any value  $\mathcal{M}$  to  $S_\phi = \{1, 2\}$  or  $S_a$ , as argued previously, we satisfy the first objective. Since  $S_a$  sends only additive shares of the maximum value, any DB owner cannot learn other DB owner’s value, except the maximum value.

## 6.4 PSI Median Query

A PSI median query over cost column corresponding to disease column over Tables 1-3 returns  $\{\langle \text{Cancer}, 300 \rangle\}$ . Note that here we first add the cost of treatment at each DB owner. However, the approach can be extended to deal with individual tuples. For solving PSI median, we extend the method of finding maximum by executing all steps as specified in §6.3 with an additional process in STEP 2. Particularly,  $S_a$  in STEP 2 of §6.3 after adding shares, sorts them, and finds the median value. If the number of DB owners is odd (even), then  $S_a$  finds the middle value (two middle values) in the sorted shares.

## 6.5 PSI Count Query

We extend PSI method (§5) to only reveal the count of common items among DB owners (i.e., the cardinality of the common item), instead of revealing common items. Recall that servers  $S_\phi$  know a permutation function  $\mathcal{PF}_{s1}$  that is not known to DB owners. The idea behind this is to find the common items over  $\chi$  and to permute the final output at servers before sending the vector (of additive share form) to DB owners. Thus, when DB owners perform computation on the vector received from servers to know the final output, the position of one in the vector will not reveal common items, while the count of one will reveal the cardinality of the common items. Thus, PSI count method follows all steps of PSI as described in §5.1 with an addition of permutation function execution by servers before sending the output to DB owners.

**PSI count information leakage discussion.** We can argue information leakage at servers and DB owners like §5.1. In addition, since both servers used the same permutation function, it will produce the correct answer at DB owner; moreover, it will hide information about the exact common items among DB owners.

**PSI Count verification.** While verifying PSI count, we cannot reveal the exact common items. Thus, in this method, we use different permutations for servers and DB owners to hide the exact position of the common item in  $\chi$ . More explicitly, each server uses two permutation

Papers	[41] & [50]	[60]	[2]	[1]	[39]	[40]	Jana [4]†	SMCQL [5]	Sharemind [7]	Conclave [65]‡	PRISM
Operations supported	PSI	PSI	PSI	PSI	PSI	PSI	PSI, PSU, aggregation	PSI via join & aggregation	PSI via join & aggregation	PSI via join & aggregation	PSI, PSU, aggregation
Verification Support	×	×	×	✓	✓	×	×	×	×	×	✓
Scalability based on experiments reported (dataset size & time)	N/A	32768 (≈50 m)	1 million (≈2 h)	32768 (≈16 m)	1 billion (≈10 m)	1000 (≈9 m)	1 million (≈1 h)	>23 million (≈23 h)	30000 (>2 h)	4 million (8 m)	20 million (At most 8 s)
Communication among servers	N/A	N/A	N/A	N/A	N/A	N/A	Yes *	Yes *	Yes *	Yes *	No
Computational Complexity	$O(n^m)$	$O(\alpha mn)$	$O(n^m)$	$O(mn^2)$	$O(mn)$ ‡‡	$O(n^m)$	$O(n^m)$	N/A *	$O(n^m)$	N/A *	$O(mX)$

**Table 12: Comparison of existing cloud-based techniques against PRISM.** Notes. (i) The scalability numbers are taken from the respective papers. (ii) Results of Sharemind [7] are taken from Conclave [65] experimental comparison. (iii) #DB owners were in each paper was reported two; thus, we executed PRISM for two DB owners for this table. (iv) Only Jana, SMCQL, Sharemind, and Conclave provide identical security like PRISM. (v) h: hours, m: minutes, s: seconds, †: We setup Jana for two DB owners each with 1M values in our experiments. ‡: Conclave [65] uses a trusted party. Yes: Requires communication among servers. No: No communication among servers. \*: Based on MPC-based systems. \*\*: N/A because executing operation in cleartext or at the trusted party. m: #DB owners. n: DB size. X: domain size. ‡‡: A insecure technique that reveals the size of the intersection, and hence fast.  $\alpha$ : The cost of Bilinear Map pairing technique.

functions  $\mathcal{PF}_{s1}$  and  $\mathcal{PF}_{s2}$ , and each DB owner utilizes two permutation functions  $\mathcal{PF}_{db1}$  and  $\mathcal{PF}_{db2}$ . Recall that these permutation function satisfy Equation 1. PSI count verification works as:

**STEP 1: DB owners.** Like PSI verification,  $\mathcal{DB}_j$  creates and sends additive shares of two tables  $\chi_j$  and  $\bar{\chi}_j$  to servers. Prior to sending shares of  $\bar{\chi}_j$ ,  $\mathcal{DB}_j$  executes  $\mathcal{PF}_{db2}$  on  $\bar{\chi}_j$ .

**STEP 2: Servers.** Servers execute Equations 6, 7 on additive shares of  $\chi_j$  and  $\bar{\chi}_j$ , implement  $\mathcal{PF}_{s1}$  on  $output^{S_\phi}[]$  and  $\mathcal{PF}_{s2}$  on  $Vout^{S_\phi}[]$ , and send permuted output vectors (each of length  $b$ ) to all DB owners.

**STEP 3: DB owners.** Finally,  $\mathcal{DB}_j$  implements  $\mathcal{PF}_{db1}$  on  $output^{S_\phi}[]$  and executes Equations 8-10 on permuted vector  $output^{S_\phi}[]$  and  $Vout^{S_\phi}[]$ . Note that due to the implementation of Equation 1 ( $\mathcal{PF}_{s1} \odot \mathcal{PF}_{db1} = \mathcal{PF}_{s2} \odot \mathcal{PF}_{db2}$ ), any  $k^{th}$  value of  $output^{S_\phi}[]$  and  $Vout^{S_\phi}[]$  will satisfy Equations 8-10 we can verify PSI count.

**Information leakage discussion.** Now we argue that why PSI count verification works and why it does not reveal the identity of the common items in the table  $\chi$ . Recall that we use four different permutations satisfying Equation 1, and  $\mathcal{PF}_i$  is unknown to servers and DB owners. Thus, neither DB owners nor servers cannot learn the exact permutation  $\mathcal{PF}_i$  by knowing their permutation functions. Hence, servers cannot link any element of  $\chi$  to  $\bar{\chi}$ , and DB owners cannot link the final output to elements of  $\chi$ . Thus,

## 6.6 Extending PSI over Multiple Attributes

In the previous sections, we explained PSI over a single attribute (or a set). We can trivially extend it to multiple attributes (or multisets). Particularly, such a query can be express in SQL as follows:

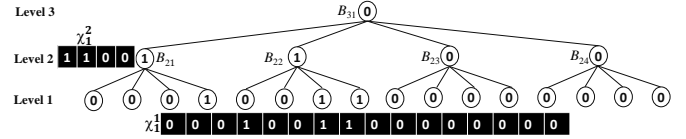
```
SELECT A_c, A_x FROM db1 INTERSECT ... INTERSECT SELECT A_c, A_x FROM db_m
```

Recall that in PSI finding method §5.1,  $\mathcal{DB}_j$  sends additive shares of a table  $\chi_j$  of length  $b = |\text{Dom}(A_c)|$ , where  $A_c$  was the attributes on which we executed PSI. Now, we can extend this method by creating a table  $\chi_j$  of length  $b = |\Pi_{i>0} \text{Dom}(A_i)|$ , where  $A_i$  are attributes on which we want to execute PSI. However, as the domain size and the number of attributes increase, such a method incurs the communication overhead. Thus, to apply the PSI method over a large (and real) domain size, as well as, to reduce the communication overhead, we provide a method, named as bucketization-based PSI.

**Optimization: bucketization-based PSI.** Before going to steps of this method, let us consider the following example:

**Example 6.6.1.** Consider two attributes  $A$  with  $|\text{Dom}(A)| = 8$  and  $B$  with  $|\text{Dom}(B)| = 2$ . Thus, DB owners can create  $\chi_j$  of 16 cells. Assume that there are two DB owners:  $\mathcal{DB}_1$  with  $\chi_1$  whose only positions 4,7,8 have one; and  $\mathcal{DB}_2$  with  $\chi_2$  whose only positions 1,6,8 have one. Thus, each DB owner sends/receives a vector of length 16 from each server.

Now, to reduce communication, we create buckets over the cell of  $\chi$  and build a tree, called *bucket-tree*, of depth  $\log_\kappa |\chi|$ , where  $\kappa$  is the number of the maximum number of child nodes that a node can have. Bucket-tree is created in a bottom-up manner by a non-overlap



**Figure 2: Bucket tree for 16 values.**

grouping of  $\kappa$  nodes, and for each level of bucket-tree a hash table (similar to  $\chi$ ) is created. Notation  $\chi_j^i$  denotes this table for  $i^{th}$  level of bucket-tree at  $\mathcal{DB}_j$ , and  $\chi_j^i[k] = 1$ , if the  $k^{th}$  node at the  $i^{th}$  level has 1.

Figure 2 shows bucket-tree for  $\mathcal{DB}_j$ ,  $|\chi| = 16$ , and  $\kappa = 4$ , with appropriate one and zero in  $\chi_j^i$ . Note that the second level shows four nodes  $B_{21}, B_{22}, B_{23}, B_{24}$  corresponding to 1–4, 5–8, 9–12, and 13–16. Since  $\mathcal{DB}_1$  has one at 4,7,8 leaf nodes, we obtain  $\chi_1^2 = \langle 1, 1, 0, 0 \rangle$ , i.e.,  $B_{21} = 1, B_{22} = 1, B_{23} = 0, B_{24} = 0$ . Here,  $B_{21} = 1$ , since its one of the child nodes has one. Now, when computing PSI,  $\mathcal{DB}_j$  can start the same computation as shown in STEP 2 of §5.1 over the specified  $i^{th}$  levels'  $\chi_j^i$ . Next, they continue the computation only for those child nodes, whose parent nodes resulted in one in STEP 3 of §5.1.

For example, in Figure 2,  $\mathcal{DB}_j$  can execute PSI for  $\chi_j^2$  and know that the only desired bucket nodes are  $B_{21}$  and  $B_{22}$  that contain common items. Thus, in the next round, they execute PSI over the first eight items of  $\chi_j^1$ , i.e., child nodes of  $B_{21}$  and  $B_{22}$ . Hence, while we use two communication rounds, DB owners/servers send  $4+8=12$  numbers instead of 16 numbers. ■

Bucketization-based PSI has the following steps:

**STEP 1A: DB owner.** Build the tree as specified in Example 6.6.1.

**STEP 1B: DB owner.** Outsource additive shares of  $i^{th}$  level's  $\chi_j^i$ .

**STEP 2: Servers.** Servers compute PSI using STEP 2 of §5.1 over  $\chi_j^i$  ( $1 \leq j \leq m$ ) and provide answers to DB owners.

**STEP 3: DB owner.**  $\mathcal{DB}_j$  computes results to find the common items in  $\chi_j^i$  and discards all non-common values of  $\chi_j^i$  and their child nodes.  $\mathcal{DB}_j$  requests servers to execute the above STEP 2 for  $\chi_j^{i-1}$  that has values corresponding to all non-discarded nodes of  $(i-1)^{th}$  level node.

**Note:** The role of DB owners in traversing the tree (i.e., the above STEP 3) can be eliminated by involving  $\mathcal{S}_a$ .

**Open problem.** In bucketization, we perform PSI at layers of the tree for eliminating ranges where corresponding child nodes have zero. However, if the data is dense (i.e., data covers most of the domain values), then bucketization-based PSI may incur overhead, since all nodes in the tree may correspond to one, leading to execute PSI on all those nodes including leaf nodes. Nevertheless, if the data is sparse (i.e., the domain is much larger than the data, as is the case of the domain to be a cartesian product of domains of two or more attributes), then higher-level nodes in the tree may have 0, leading to eliminate ranges of domain on which PSI is performed. Developing an optimal bucketization strategy that minimizes PSI execution is an interesting open problem.



## 7 PRIVATE SET UNION (PSU) QUERY

This section develops a method for finding union (denoted by PSU) among  $m > 1$  different DB owners over an attribute  $A_c$  (which is assumed to exist at all DB owners).

**High-level idea.** Likewise PSI method (as presented in §5), each DB owner uses a publicly known hash function to map distinct values of  $A_c$  attribute in a table of cells at most  $|\text{Dom}(A_c)|$ , where  $|\text{Dom}(A_c)|$  refers to the size of the domain of  $A_c$ , and outsources each element of the table in additive share form to *two servers*  $\mathcal{S}_\phi$ ,  $\phi \in \{1, 2\}$ .  $\mathcal{S}_\phi$  computes the union obliviously, thereby DB owners will receive a vector of length  $|\text{Dom}(A_c)|$  having either 0 or 1 of additive shared form. After adding the share for an  $i^{\text{th}}$  element, DB owners only know whether the element is in the union or not; nothing else.

**STEP 1: DB owner.** This step is identical to STEP 1 of PSI (see §5.1).

**STEP 2: Server.** Each server  $\mathcal{S}_\phi$  ( $\phi \in \{1, 2\}$ ) holds the  $\phi^{\text{th}}$  additive share of the table  $\chi$  of  $m$  DB owners and executes the following operation:

$$\text{rand}[] \leftarrow \mathcal{PRG}(\text{seed}) \quad (20)$$

Each server  $\mathcal{S}_\phi$  performs the following operations: (i) generates  $b$  pseudorandom numbers, (ii) performs (arithmetic) addition of the  $i^{\text{th}}$  additive secret-shares from all DB owners, (iii) multiplies the resultant of the previous step with  $i^{\text{th}}$  pseudorandom number and then takes modulo, and (iv) sends  $b$  results to all DB owners.

**STEP 3: DB owner.** On receiving two vectors, each of length  $b$ , from two servers, DB owners execute modular addition over  $i^{\text{th}}$  shares of both vectors to know the final answer (Equation 21). It results in either zero or any random number, where zero shows that the  $i^{\text{th}}$  element of  $\chi$  is not present at any DB owner, while a random number shows the  $i^{\text{th}}$  element of  $\chi$  is present at one of the DB owners.

$$\text{fop}_i \leftarrow (\text{output}_i^{S_1} + \text{output}_i^{S_2}) \bmod \delta \quad (21)$$

**Correctness.** When we plug Equation 20 into Equation 21, we obtain:

$$\begin{aligned} \text{fop}_i &= (((\sum_{j=1}^m A(x_i)_j^1) \times \text{rand}[i]) \bmod \delta) \\ &+ (((\sum_{j=1}^m A(x_i)_j^2) \times \text{rand}[i]) \bmod \delta) \bmod \delta \\ &= (((\sum_{j=1}^m A(x_i)_j^1) \times \text{rand}[i]) + ((\sum_{j=1}^m A(x_i)_j^2) \times \text{rand}[i])) \bmod \delta \\ &= ((\sum_{j=1}^m A(x_i)_j) \times \text{rand}[i]) \bmod \delta \end{aligned}$$

Thus, whenever the  $i^{\text{th}}$  element of  $\chi$  will present at any DB owner, it will result in a random number; otherwise, zero due to  $\sum_{j=1}^m (x_i)_j = 0$ . Further, since servers generated random numbers between 1 to  $\delta - 1$ , any  $\text{fop}_i$  that should be a random number, will never be zero at DB owners, due to not using a random number  $\delta$ , i.e.,  $((\sum_{j=1}^m A(x_i)_j) \times \delta) \bmod \delta$ .

**Information leakage discussion.** We need to prevent information leakage at servers and at DB owners.

(1) *Server perspective.* The servers only know  $\delta$ . However, based on  $\delta$ , an individual server cannot reconstruct  $\chi$  that is sent by DB owners in additive shared form. Servers execute an identical operation on all shares of  $m$  DB owners; thus, access-patterns are hidden from servers preventing them to know anything based on access-patterns. Since each output contains an identical number of bits, it does not reveal to the adversary based on the output size.

(2) *DB owner perspective.* We need to hide from all DB owners the fact that how many DB owners do not have one at the  $i^{\text{th}}$  position of  $\chi$ . Revealing this can reveal the intersection of all elements in

$\chi$  and provide additional information to DB owners. Note that DB owners can learn such information, if they receive the following:  $\text{output}_i^{S_\phi} \leftarrow (\sum_{j=1}^m A(x_i)_j^\phi) \bmod \delta$ . Since servers multiply the result of this by a random number, DB owner cannot learn the actual number of ones at the  $i^{\text{th}}$  position of  $\chi$ , unless they know the random number.

## 8 EXPERIMENTAL EVALUATION

This section evaluates the scalability of PRISM on different-sized datasets and a different number of DB owners. Also, we evaluate the verification overhead and compare it against other MPC-based systems. We used a 16GB RAM machine with 4 cores for each of the DB owners and [three AWS servers of 32GB RAM, 3.5GHz Intel Xeon CPU with 16 cores to store shares](#). The communication between DB owners and servers were done using the `scp` protocol, and  $\eta, \delta$  were 227, 113, respectively.

### 8.1 PRISM Evaluation

**Dataset generation.** We used five columns (Orderkey (OK), Partkey (PK), Linenummer (LN), Suppkey(SK), and Discount (DT)) of LineItem table of TPC-H benchmark. We experimented with domain sizes (i.e., the number of values) of 5M and 20M for the *OK column* on which we executed PSI and PSU. Further, we selected at most **50 DB owners**. To the best of our knowledge, this is the first such experiment of multi-owner large datasets. OK column is used for PSI/PSU, and other columns were used for aggregation operations. To generate secret-shared dataset, each DB owner maintained a LineItem table containing at most 5M (20M) OK values. To outsource the database, each DB owner did the following:

Created a table of 11 columns, as shown in Table 13, in which the first five columns contain the secret-shared data of LineItem table, the next five columns contain the verification data for the first five columns, and the last column (aOK) was used for computing the average. All verification column names are prefixed with the character ‘v.’

Real data column					For verification					Average
OK	PK	LN	SK	DT	vOK	vPK	vLN	vSK	vDT	aOK

**Table 13: Table structure created by PRISM.**

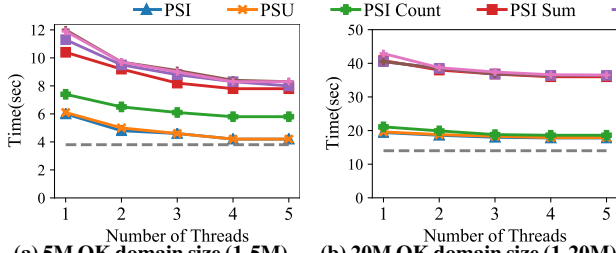
First column of Table 13 was created over OK column of LineItem table (by following STEP 1 of §5.1) for executing PSI/PSU over OK. vOK column was created to verify PSI results (by following STEP 1 of §5.2). Columns PK and vPK were created using the following command: `select OK, sum(PK) from LineItem group by OK`. Other columns (LN, SK, DT, vLN, vSK, vDT) were created by using similar SQL commands.

Columns aOK was created using the following command: `select count(*) from LineItem group by OK`.

Finally, permute all values of verification columns and create additive shares of (OK and vOK), as well as, multiplicative shares of all remaining columns.

**Share generation time.** The time to generate two additive shares and three multiplicative shares of the respective first five columns of Table 13 in the case of 5M (or 20M) OK domain size was 121s (or 548s). Furthermore, the time for creating each additional column for verification took 20s (or 90s) in the case of 5M (or 20M) domain values.

**Exp 1. PRISM performance on multi-threaded implementation at AWS.** Since identical computations are executed on each row of the table, we exploit multiple CPU cores by writing a parallel implementation of PRISM. The parallel implementation divides rows into multiple blocks with each thread processing a single block. We



(a) 5M OK domain size (1-5M). (b) 20M OK domain size (1-20M).  
Figure 3: Exp 1. PRISM performance on multi-threaded implementation at AWS.

Data size	Sum over different attributes				Max over different attributes			
	1	2	3	4	1	2	3	4
5M	8.2	12.1	15.9	20.4	10	14.6	19	23.5
20M	33.4	48.6	63.5	81.9	36.6	53.3	70	87.4

Table 14: Exp 1. Multi-column aggregation query performance (time in seconds).

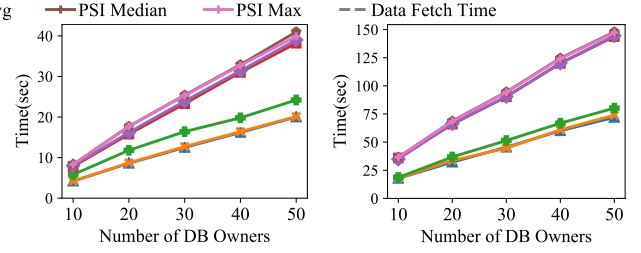
increased the number of threads from 1 to 5; see Figure 3, while fixing DB owners to 10. Increasing threads more than 5 did not provide speed-up, since reading/writing of data quickly becomes the bottleneck as the number of threads increase. Observe that the data fetch time from the database remains (almost) identical; see Figure 3. *PSI and PSU queries.* Figure 3 shows the time taken by PSI/PSU over the OK column. Observe that as the number of values in OK column increases (from 5M to 20M), the time increases (almost) linearly from 4s to 18s, respectively.

*Aggregation queries over PSI.* We executed PSI count, average, sum, maximum, and median queries; see Figure 3. Observe that the processing time of PSI count is almost the same as that of PSI, since it involves only one round of computation in which we permute the output of PSI. In contrast, other aggregation operations (sum, average, maximum, and median) incur almost twice processing cost at servers, since they involve computing PSI over OK column in the first round and, then, computing aggregation in the second round. For this experiment, we computed sum only over DT column and maximum/median over PK column. Table 14 shows the impact of computing sum and maximum over multiple attributes (from 1 to 4). As we increase the number of attributes, the computation of respective aggregation operation also increases, due to additional addition/multiplication/modulo operations on additional attributes.

**Exp 2. Impact of the number of DB owners.** Since we developed PRISM to deal with multiple DB owners, we investigated the impact of DB owners by selecting 10, 20, 30, 40, 50 DB owners, for two different domain sizes of OK column. Figure 4 shows the server processing time for PSI, PSU, and aggregation over PSI. Note that as the number of DB owners increases, the computation time at the server increases linearly, due to linearly increasing number of addition/multiplication/modulo operations; e.g., on 5M OK values, PSI processing took 4.2s, 8.6s, 12.5s, 16.2s, and 20s in the case of 10, 20, 30, 40, 50 DB owners.

**Exp 3. Result verification overheads.** Figure 5 shows the overheads of the result verification approaches on 5M and 20M domain values of OK column and 10 DB owners. *PSI verification and PSI count verification* took almost twice processing time than respective non-verification methods, due to executing additional operations on the same amount of data for verification. *PSI sum verification* took ( $\approx 20$ s on 5M) more than two times than non-verification based PSI sum ( $\approx 7$ s on 5M), since we verified both PSI and sum.

**Exp 4. DB owner processing time in result construction.** PRISM requires DB owners to perform computation on additive or multiplicative shares. Table 15 shows the processing time at a DB owner over 5M and 20M domain values for different operations. It is clear that the DB owner processing time is significantly less than the server



(a) 5M OK domain size (1-5M). (b) 20M OK domain size (1-20M).  
Figure 4: Exp 2. PRISM dealing with multiple DB owners.

Data Size	PSI	Count	Sum	Avg	Max	PSU	vPSI	vCount	vSum
5M	1.3	1.7	3.1	3.2	2.8	1.3	2.8	1.7	4.1
20M	4.8	5.4	10.3	10.3	9.5	4.8	11.6	5.6	13.2

Table 15: Exp 4. DB owner processing time in result construction (time in seconds).

processing time. In case of 5M (20M) OK values and 50 DB owners, each DB owner took at most 4s (13s) in PSI Sum (PSI Sum) query, while servers took at least 20s (72s) in PSI (Sum) query; see Figure 4.

**Exp 5. Impact of communication cost.** PRISM protocols involve at most two rounds, where servers send data of size equal to the domain size in the first and second rounds of query execution. Thus, it is required to measure the impact of communication cost, since it may affect the overall performance. Among the proposed protocols, the maximum amount of data flows for maximum/median queries, due to first receiving the answers of PSI, then additive share transmission from each DB owner to a server, and finally, receiving the answer of the maximum query from a server to DB owners. Here, the overall data was transmitted of size 60MB (240MB) in the case of 5M (20M) OK values and took 1.2s (4.8s), 0.6s (2.4s), 0.1s (0.4s) on slow (50MB/s), medium (100MB/s), and fast (500MB/s) speed of data transmission. *To measure the communication cost, we simulated network cost by finding appropriate delays in the transmission, where delay was determined by dividing data size by the network speed.*

**Exp 6. Impact of bucketization.** Figure 6 shows the reduction in the number of values on which we need to execute PSI when using bucketization technique (explained in §6.6). For our experiment, we created a tree with fanout of 10, height 9 and 100M values at the leaf level. In Figure 6, we refer to the percentage of leaf nodes of the tree that containing one as *fill factor*. We use a term *actual domain size* (in Figure 6) that refers to the number of items on which we execute PSI. Note that actual domain size is different from *real domain size* that refers to the domain values given to us, i.e., 100M. Note that the actual domain size depends on the fill factor and impacts the performance of PSI. Observe that when the fill factor is 100% (i.e., all leaf nodes have one, and thus, the entire tree has one), the actual domain size was 111M. In contrast, if the fill factor was only 0.01% of 100M values (i.e., 10K), then most of the tree contained zero; thus, we run PSI only on actual domain size equal to 400K, instead of real domain size of 100M. Note that for this experiment, we generated the data randomly. If there is a correlation in the data (which is the case in most real-world datasets), bucketization results will be even better.

## 8.2 Comparing with Other Works

We compare PRISM against the state-of-the-art cloud-based industrial MPC-based systems: Galois Inc.'s Jana [4], since it provides the identical security guarantees at servers as PRISM. To evaluate Jana, we inserted two LineItem tables (each of 1M rows) having (OK, PK, LN, SK, DT) columns and executed join on OK column. However, the join execution took more than 1 hour to complete.

[1, 2, 39–41, 50, 60] provide cloud-based PSI/PSU/aggregation techniques/systems. *We could not experimentally compare Prism*

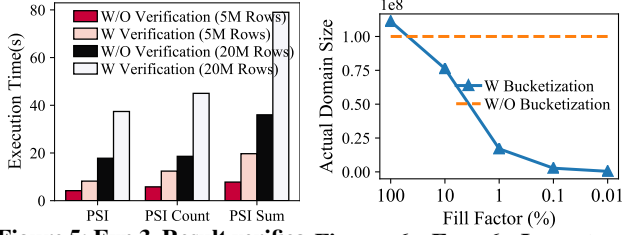


Figure 5: Exp 3. Result verification overheads. Figure 6: Exp 6. Impact of bucketization.

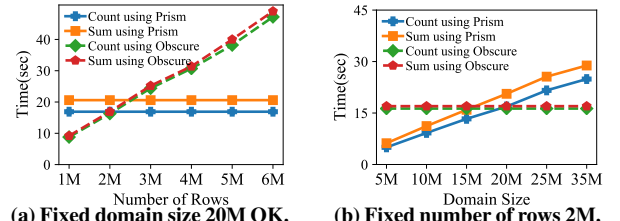
against such systems, since none of them is not open source.<sup>4</sup> Thus, in Table 12, we report experimental results from those papers, just for intuition purposes. With the exception of [39], none of the techniques supports large-sized dataset, has quadratic/exponential complexity or uses expensive cryptographic techniques [60]. While [39] scales better, it does not support aggregation and, moreover, reveals which item is in the intersection set. For a fair comparison, we report PRISM results only for two DB owners in Table 12, since other papers do not provide experimental results for more than two DB owners. Recall that in our experiments (Figure 4a), PRISM supports 50 DB owners and takes at most  $\approx 41$  seconds on 5M values. Further note that, in the case of 1B values and two DB owners, PRISM takes  $\approx 7.3$  mins, unlike [39] that took  $\approx 10$  mins.

There are several non-cloud-based PSI approaches. However, such approaches cannot be directly compared against PRISM, due to a different model used (in which DB owners communicate amongst themselves and do not outsource data to the cloud) and/or different security properties. Just to put some numbers in this context, recent work [45] took 304s in the case of 14 DB owners each with 1M values, and [36] took at least  $\approx 400$ s for PSI sum on 100K values.

**Other PSI/PSU and max/min finding protocols.** Several PSI protocols have been proposed [1, 3, 14, 18, 19, 21, 23, 24, 26, 28, 33, 34, 39, 41, 42, 44–48, 52, 53, 55, 56, 58, 59, 61, 64, 69], and a survey of PSI protocols may be found in [55]. Among these techniques, only [1, 2, 4, 5, 7, 39–41, 50, 60, 65] are developed for the cloud settings; as we compared in Table 12. Since the classic Millionaire’s problem has been proposed by Yao [68], many schemes for comparison/maximum finding were proposed; e.g., [8–10, 20, 32, 54, 63]. However, such techniques show limitations: many communication rounds, restricted to two DB owners, quadratic computation cost at servers, not dealing with malicious adversaries in the cloud setting, and/or no support for result verification.

**Comparison between PRISM and OBSCURE.** While both PRISM and OBSCURE [30] are based on secret-sharing, they are significantly different from each other in terms of: (i) purposes: PRISM is for computing PSI/PSU queries over multi-owner databases, while OBSCURE is for query processing over outsourced data and does not support PSI/PSU queries; (ii) implementation of secret-sharing: PRISM is based on domain-based representation, while OBSCURE is based on unary representation; (iii) offered functionalities: PRISM provides aggregation over PSI/PSI, while OBSCURE provides complex conjunctive and disjunctive aggregation queries; and (iv) query execution complexities: PRISM complexity is upper bounded by  $m \times \text{Dom}(A_c)$ , where  $m$  is the number of DB owners and  $\text{Dom}(A_c)$  is the domain of the attribute  $A_c$ , while OBSCURE complexity is upper bounded by  $n \times L$ , where  $n$  is the number of tuples and  $L$  is the length of a value in unary representation. Thus, a direct comparison between the two non-identical systems is infeasible. Nonetheless, it is interesting to see the overheads of the different secret-sharing

<sup>4</sup>None of these techniques have open sources implementations, except [5]. We installed [5] that works for a very small data and incurs runtime errors. We have reported this issue to the author as well.



(a) Fixed domain size 20M OK. (b) Fixed number of rows 2M. Figure 7: PRISM vs OBSCURE comparison.

techniques. We can do this by mimicking aggregation queries with a simple selection that OBSCURE supports, as a PSI query.

For a practical comparison of the two systems, we took two attributes OK and PK of LineItem table and outsourced them by following OBSCURE and PRISM share generation process. We executed one-dimensional count and sum queries using the methods of PRISM and OBSCURE. In the first experiment (using a single thread), we fixed the domain size of OK column to 20M and increased the number of rows as shown on x-axis of Figure 7a. Figure 7a shows that executing count/sum query using PRISM takes an identical time, regardless of dataset size, due to PRISM’s dependence on the domain size. The same argument holds for sum query under PRISM. However, in the case of OBSCURE, the computation time for both queries increases as the dataset increases. In the second experiment Figure 7b, we fixed the number of rows to be 2M, while varying the domain size. Note that in this case, as the domain size increases, PRISM computation time increases, while OBSCURE computation time remains the same.

## 9 CONCLUSION

This paper describes PRISM based on secret-sharing that allows multiple DB owners to outsource data to (a majority of) non-colluding servers that can behave like honest-but-curious servers and malicious servers in terms of the computation that they perform. It exploits the additive and multiplicative homomorphic property of secret-sharing techniques to implement both set operations and aggregation functions efficiently. Experimental results show PRISM scales to both a large number of DB owners and to large datasets, compared to existing systems. Future directions include dealing with: (i) multiple attributes more efficiently than bucketization, (ii) dealing with malicious DB owners, and (iii) a broader set of SQL queries.



## REFERENCES

- [1] A. Abadi et al. VD-PSI: verifiable delegated private set intersection on outsourced private datasets. In *FC*, pages 149–168, 2016.
- [2] A. Abadi et al. Efficient delegated private set intersection on outsourced private datasets. *IEEE Trans. Dependable Secur. Comput.*, 16(4):608–624, 2019.
- [3] T. Araki et al. High-throughput semi-honest secure three-party computation with an honest majority. In *CCS*, pages 805–817, 2016.
- [4] D. W. Archer et al. From keys to databases - real-world applications of secure multi-party computation. *Comput. J.*, 61(12):1749–1771, 2018.
- [5] J. Bader et al. SMCQL: secure query processing for private data networks. *Proc. VLDB Endow.*, 10(6):673–684, 2017.
- [6] M. Blum et al. How to generate cryptographically strong sequences of pseudo-random bits. *SIAM J. Comput.*, 13(4):850–864, 1984.
- [7] D. Bogdanov et al. Sharemind: A framework for fast privacy-preserving computations. In *ESORICS*, pages 192–206, 2008.
- [8] D. Bogdanov et al. A practical analysis of oblivious sorting algorithms for secure multi-party computation. In *NordSec*, pages 59–74, 2014.
- [9] M. Burkhart et al. Fast privacy-preserving top-k queries using secret sharing. In *ICCCN*, pages 1–7, 2010.
- [10] M. Burkhart et al. SEPIA: privacy-preserving aggregation of multi-domain network events and statistics. In *USENIX Security Symposium*, pages 223–240, 2010.
- [11] R. Canetti. Security and composition of multiparty cryptographic protocols. *J. Cryptology*, 13(1):143–202, 2000.
- [12] R. Canetti et al. Adaptively secure multi-party computation. In G. L. Miller, editor, *STOC*, pages 639–648, 1996.
- [13] D. Cash et al. Leakage-abuse attacks against searchable encryption. In *CCS*, pages 668–679, 2015.
- [14] H. Chen et al. Fast private set intersection from homomorphic encryption. In *CCS*, pages 1243–1255, 2017.
- [15] J. H. Cheon et al. Multi-party privacy-preserving set intersection with quasi-linear complexity. *IEICE Trans. Fundam. Electron. Commun. Comput. Sci.*, 95-A(8):1366–1378, 2012.
- [16] K. Chida et al. An efficient secure three-party sorting protocol with an honest majority. *IACR Cryptol. ePrint Arch.*, 2019:695, 2019.
- [17] R. M. Corless and N. Fillion. A graduate introduction to numerical methods. *AMC*, 10:12, 2013.
- [18] E. D. Cristofaro et al. Linear-complexity private set intersection protocols secure in malicious model. In *ASIACRYPT*, pages 213–231, 2010.
- [19] E. D. Cristofaro et al. Fast and private computation of cardinality of set intersection and union. In *CANS*, pages 218–231, 2012.
- [20] I. Damgård et al. Unconditionally secure constant-rounds multi-party computation for equality, comparison, bits and exponentiation. In *TCC*, pages 285–304, 2006.
- [21] C. Dong et al. When private set intersection meets big data: an efficient and scalable protocol. In *CCS*, pages 789–800, 2013.
- [22] R. Egert et al. Privately computing set-union and set-intersection cardinality via bloom filters. In *ACISP*, pages 413–430, 2015.
- [23] B. H. Falk et al. Private set intersection with linear communication from general assumptions. In *WPES@CCS*, pages 14–25, 2019.
- [24] M. J. Freedman et al. Efficient private matching and set intersection. In *EUROCRYPT*, pages 1–19, 2004.
- [25] M. J. Freedman et al. Keyword search and oblivious pseudorandom functions. In *TCC*, pages 303–324, 2005.
- [26] J. Furukawa et al. High-throughput secure three-party computation for malicious adversaries and an honest majority. In *EUROCRYPT*, pages 225–255, 2017.
- [27] O. Goldreich et al. How to construct random functions. *J. ACM*, 33(4):792–807, 1986.
- [28] O. Goldreich et al. How to play any mental game or A completeness theorem for protocols with honest majority. In *STOC*, pages 218–229, 1987.
- [29] D. M. Goldschlag et al. Onion routing. *Commun. ACM*, 42(2):39–41, 1999.
- [30] P. Gupta et al. Obscure: Information-theoretic oblivious and verifiable aggregation queries. *Proc. VLDB Endow.*, 12(9):1030–1043, 2019.
- [31] H. Hacigümüş et al. Executing SQL over encrypted data in the database-service-provider model. In *SIGMOD*, pages 216–227, 2002.
- [32] K. Hamada et al. Practically efficient multi-party sorting protocols from comparison sort algorithms. In *ICISC*, pages 202–216, 2012.
- [33] C. Hazay et al. Scalable multi-party private set-intersection. In *PKC*, pages 175–203, 2017.
- [34] Y. Huang et al. Private set intersection: Are garbled circuits better than custom protocols? In *NDSS*, 2012.
- [35] R. Inbar et al. Efficient scalable multiparty private set-intersection via garbled bloom filters. In *SCN*, pages 235–252, 2018.
- [36] M. Ion et al. On deploying secure computing commercially: Private intersection-sum protocols and their business applications. *IACR Cryptol. ePrint Arch.*, 2019:723, 2019.
- [37] M. S. Islam et al. Access pattern disclosure on searchable encryption: Ramification, attack and mitigation. In *NDSS*, 2012.
- [38] W. Jiang et al. Transforming semi-honest protocols to ensure accountability. *Data Knowl. Eng.*, 65(1):57–74, 2008.
- [39] S. Kamara et al. Scaling private set intersection to billion-element sets. In *FC*, pages 195–215, 2014.
- [40] F. Kerschbaum. Collusion-resistant outsourcing of private set intersection. In *SAC*, pages 1451–1456, 2012.
- [41] F. Kerschbaum. Outsourced private set intersection using homomorphic encryption. In *ASIACCS*, pages 85–86, 2012.
- [42] L. Kissner et al. Privacy-preserving set operations. In *CRYPTO*, pages 241–257, 2005.
- [43] V. Kolesnikov et al. Efficient batched oblivious PRF with applications to private set intersection. *IACR Cryptol. ePrint Arch.*, 2016:799, 2016.
- [44] V. Kolesnikov et al. Efficient batched oblivious PRF with applications to private set intersection. In *CCS*, pages 818–829, 2016.
- [45] V. Kolesnikov et al. Practical multi-party private set intersection from symmetric-key techniques. In *CCS*, pages 1257–1272, 2017.
- [46] P. H. Le et al. Two-party private set intersection with an untrusted third party. In *CCS*, pages 2403–2420, 2019.
- [47] R. Li et al. An unconditionally secure protocol for multi-party set intersection. In *Applied Cryptography and Network Security*, pages 226–236, 2007.
- [48] Y. Li et al. Delegatable order-revealing encryption. In *AsiaCCS*, pages 134–147, 2019.
- [49] Y. Lindell. Secure multiparty computation (MPC). *IACR Cryptol. ePrint Arch.*, 2020:300, 2020.
- [50] F. Liu et al. Encrypted set intersection protocol for outsourced datasets. In *ICCE*, pages 135–140, 2014.
- [51] S. Madden et al. TAG: A tiny aggregation service for ad-hoc sensor networks. In *OSDI*, 2002.
- [52] D. Many et al. Fast private set operations with sepi. *ETZ G93*, 2012.
- [53] G. S. Narayanan et al. Multi party distributed private matching, set disjointness and cardinality of set intersection with information theoretic security. In *CANS*, pages 21–40, 2009.
- [54] T. Nishide et al. Multiparty computation for interval, equality, and comparison without bit-decomposition protocol. In *PKC*, pages 343–360, 2007.
- [55] B. Pinkas et al. Faster private set intersection based on OT extension. In *USENIX Security*, pages 797–812, 2014.
- [56] B. Pinkas et al. Phasing: Private set intersection using permutation-based hashing. In *USENIX Security*, pages 515–530, 2015.
- [57] B. Pinkas et al. Scalable private set intersection based on OT extension. *ACM Trans. Priv. Secur.*, 21(2):7:1–7:35, 2018.
- [58] B. Pinkas et al. SpOT-Light: Lightweight private set intersection from sparse OT extension. In *CRYPTO*, pages 401–431, 2019.
- [59] B. Pinkas et al. PSI from paxos: Fast, malicious private set intersection. In *EUROCRYPT*, pages 739–767, 2020.
- [60] S. Qiu et al. Identity-based private matching over outsourced encrypted datasets. *IEEE Trans. Cloud Comput.*, 6(3):747–759, 2018.
- [61] P. Rindal et al. Malicious-secure private set intersection via dual execution. In *CCS*, pages 1229–1242, 2017.
- [62] A. Shamir. How to share a secret. *Communication of ACM*, 22(11):612–613, 1979.
- [63] J. Vaidya et al. Privacy-preserving top-k queries. In *ICDE*, pages 545–546, 2005.
- [64] J. Vaidya et al. Secure set intersection cardinality with application to association rule mining. *J. Comput. Secur.*, 13(4):593–622, 2005.
- [65] N. Volgushev et al. Conclave: secure multi-party computation on big data. In *EuroSys*, pages 3:1–3:18, 2019.
- [66] C. Wang et al. Secure ranked keyword search over encrypted cloud data. In *ICDCS*, pages 253–262, 2010.
- [67] F. Wang et al. Splinter: Practical private queries on public data. In *NSDI*, pages 299–313, 2017.
- [68] A. C. Yao. How to generate and exchange secrets (extended abstract). In *FOCS*, pages 162–167, 1986.
- [69] E. Zhang et al. Efficient multi-party private set intersection against malicious adversaries. In *CCSW*, page 93–104, 2019.
- [70] Q. Zheng et al. Verifiable delegated set intersection operations on outsourced encrypted data. In *IC2E*, pages 175–184, 2015.